

This Page Is Inserted by IFW Operations  
and is not a part of the Official Record

## **BEST AVAILABLE IMAGES**

Defective images within this document are accurate representations of the original documents submitted by the applicant.

Defects in the images may include (but are not limited to):

- BLACK BORDERS
- TEXT CUT OFF AT TOP, BOTTOM OR SIDES
- FADED TEXT
- ILLEGIBLE TEXT
- SKEWED/SLANTED IMAGES
- COLORED PHOTOS
- BLACK OR VERY BLACK AND WHITE DARK PHOTOS
- GRAY SCALE DOCUMENTS

**IMAGES ARE BEST AVAILABLE COPY.**

**As rescanning documents *will not* correct images,  
please do not report the images to the  
Image Problem Mailbox.**





(12) **EUROPEAN PATENT SPECIFICATION**

- (45) Date of publication of patent specification: **22.06.94** (51) Int. Cl.<sup>5</sup>: **C12N 15/10, C12P 21/02, G01N 33/534, C12Q 1/68**  
 (21) Application number: **90124234.7**  
 (22) Date of filing: **25.04.83**  
 (60) Publication number of the earlier application in accordance with Art.76 EPC: **0 108 128**

(54) **The manufacture and expression of large structural genes.**

(30) Priority: **06.05.82 US 375494**  
**15.04.83 US 483451**

(43) Date of publication of application:  
**02.05.91 Bulletin 91/18**

(45) Publication of the grant of the patent:  
**22.06.94 Bulletin 94/25**

(84) Designated Contracting States:  
**AT BE CH DE FR GB LI LU NL SE**

(56) References cited:

**NUCLEIC ACIDS RESEARCH, Symposium Series, no. 11, 24th November 1982, pages 29-32, IRL Press Ltd, Oxford, GB; S. TANAKA et al.: "Expression in Escherichiacoli of chemically synthesized gene for a human immune interferon"**

**Nature, vol. 292 (1981), pages 756-762**

(73) Proprietor: **AMGEN INC.**  
**1840 Dehavilland Drive**  
**Thousand Oaks California 91320 -1789(US)**

(72) Inventor: **Alton, Norman K.**  
**3320 Silver Spur Court**  
**Thousand Oaks, California 91360(US)**  
 Inventor: **Peters, Mary A.**  
**945, 37th Street**  
**Boulder, CO 80303(US)**  
 Inventor: **Stabinsky, Yitzhak**  
**3415 Heidelberg Drive**  
**Boulder, CO 80303(US)**  
 Inventor: **Snitman, David L.**  
**1475 Ithica Drive**  
**Boulder, CO 80303(US)**

(74) Representative: **Brown, John David et al**  
**FORRESTER & BOEHMERT**  
**Franz-Joseph-Strasse 38**  
**D-80801 München (DE)**

Note: Within nine months from the publication of the mention of the grant of the European patent, any person may give notice to the European Patent Office of opposition to the European patent granted. Notice of opposition shall be filed in a written reasoned statement. It shall not be deemed to have been filed until the opposition fee has been paid (Art. 99(1) European patent convention).

**EP 0 424 990 B1**

## Description

The present invention relates generally to the manipulation of genetic materials and, more particularly, to the manufacture of specific DNA sequences useful in recombinant procedures to secure the production of proteins of interest.

Genetic materials may be broadly defined as those chemical substances which program for and guide the manufacture of constituents of cells and viruses and direct the responses of cells and viruses. A long chain polymeric substance known as deoxyribonucleic acid (DNA) comprises the genetic material of all living cells and viruses except for certain viruses which are programmed by ribonucleic acids (RNA). The repeating units in DNA polymers are four different nucleotides, each of which consists of either a purine (adenine or guanine) or a pyrimidine (thymine or cytosine) bound to a deoxyribose sugar to which a phosphate group is attached. Attachment of nucleotides in linear polymeric form is by means of fusion of the 5' phosphate of one nucleotide to the 3' hydroxyl group of another. Functional DNA occurs in the form of stable double stranded associations of single strands of nucleotides (known as deoxyligonucleotides), which associations occur by means of hydrogen bonding between purine and pyrimidine bases [i.e., "complementary" associations existing either between adenine (A) and thymine (T) or guanine (G) and cytosine (C)]. By convention, nucleotides are referred to by the names of their constituent purine or pyrimidine bases, and the complementary associations of nucleotides in double stranded DNA (i.e., A-T and G-C) are referred to as "base pairs". Ribonucleic acid is a polynucleotide comprising adenine, guanine, cytosine and uracil (U), rather than thymine, bound to ribose and a phosphate group.

Most briefly put, the programming function of DNA is generally effected through a process wherein specific DNA nucleotide sequences (genes) are "transcribed" into relatively unstable messenger RNA (mRNA) polymers. The mRNA, in turn, serves as a template for the formation of structural, regulatory and catalytic proteins from amino acids. This translation process involves the operations of small RNA strands (tRNA) which transport and align individual amino acids along the mRNA strand to allow for formation of polypeptides in proper amino acid sequences. The mRNA "message", derived from DNA and providing the basis for the tRNA supply and orientation of any given one of the twenty amino acids for polypeptide "expression", is in the form of triplet "codons" -- sequential groupings of three nucleotide bases. In one sense, the formation of a protein is the ultimate form of "expression" of the programmed genetic message provided by the nucleotide sequence of a gene.

Certain DNA sequences which usually "precede" a gene in a DNA polymer provide a site for initiation of the transcription into mRNA. These are referred to as "promoter" sequences. Other DNA sequences, also usually "upstream" of (i.e., preceding) a gene in a given DNA polymer, bind proteins that determine the frequency (or rate) of transcription initiation. These other sequences are referred to as "regulator" sequences. Thus, sequences which precede a selected gene (or series of genes) in a functional DNA polymer and which operate to determine whether the transcription (and eventual expression) of a gene will take place are collectively referred to as "promoter/regulator" or "control" DNA sequences. DNA sequences which "follow" a gene in a DNA polymer and provide a signal for termination of the transcription into mRNA are referred to as "terminator" sequences.

A focus of microbiological processing for nearly the last decade has been the attempt to manufacture industrially and pharmaceutically significant substances using organisms which do not initially have genetically coded information concerning the desired product included in their DNA. Simply put, a gene that specifies the structure of a product is either isolated from a "donor" organism or chemically synthesized and then stably introduced into another organism which is preferably a self-replicating unicellular microorganism. Once this is done, the existing machinery for gene expression in the "transformed" host cells operates to construct the desired product.

The art is rich in patent and literature publications relating to "recombinant DNA" methodologies for the isolation, synthesis, purification and amplification of genetic materials for use in the transformation of selected host organisms. U.S. Letters Patent No. 4,237,224 to Cohen, et al., for example, relates to transformation of procaryotic unicellular host organisms with "hybrid" viral or circular plasmid DNA which includes selected exogenous DNA sequences. The procedures of the Cohen, et al. patent first involve manufacture of a transformation vector by enzymatically cleaving viral or circular plasmid DNA to form linear DNA strands. Selected foreign DNA strands are also prepared in linear form through use of similar enzymes. The linear viral or plasmid DNA is incubated with the foreign DNA in the presence of ligating enzymes capable of effecting a restoration process and "hybrid" vectors are formed which include the selected foreign DNA segment "spliced" into the viral or circular DNA plasmid.

Transformation of compatible unicellular host organisms with the hybrid vector results in the formation of multiple copies of the foreign DNA in the host cell population. In some instances, the desired result is

simply the amplification of the foreign DNA and the "product" harvested is DNA. More frequently, the goal of transformation is the expression by the host cells of the foreign DNA in the form of large scale synthesis of isolatable quantities of commercially significant protein or polypeptide fragments coded for by the foreign DNA. See also, e.g., U.S. Letters Patent Nos. 4,269,731 (to Shine), 4,273,875 (to Manis) and 4,293,652 (to Cohen).

The success of procedures such as described in the Cohen, et al. patent is due in large part to the ready availability of "restriction endonuclease" enzymes which facilitate the site-specific cleavage of both the unhybridized DNA vector and, e.g., eukaryotic DNA strands containing the foreign sequences of interest. Cleavage in a manner providing for the formation of single stranded complementary "ends" on the double stranded linear DNA strands greatly enhances the likelihood of functional incorporation of the foreign DNA into the vector upon "ligating" enzyme treatment. A large number of such restriction endonuclease enzymes are currently commercially available [See, e.g., "BRL Restriction Endonuclease Reference Chart" appearing in the "'81/'82 Catalog" of Bethesda Research Laboratories, Inc., Gaithersburg, Maryland.] Verification of hybrid formation is facilitated by chromatographic techniques which can, for example, distinguish the hybrid plasmids from non-hybrids on the basis of molecular weight. Other useful verification techniques involve radioactive DNA hybridization.

Another manipulative "tool" largely responsible for successes in transformation of procaryotic cells is the use of selectable "marker" gene sequences. Briefly put, hybrid vectors are employed which contain, in addition to the desired foreign DNA, one or more DNA sequences which code for expression of a phenotypic trait capable of distinguishing transformed from non-transformed host cells. Typical marker gene sequences are those which allow a transformed procaryotic cell to survive and propagate in a culture medium containing metals, antibiotics, and like components which would kill or severely inhibit propagation of non-transformed host cells.

Successful expression of an exogenous gene in a transformed host microorganism depends to a great extent on incorporation of the gene into a transformation vector with a suitable promoter/regulator region present to insure transcription of the gene into mRNA and other signals which insure translation of the mRNA message into protein (e.g., ribosome binding sites). It is not often the case that the "original" promoter/regulator region of a gene will allow for high levels of expression in the new host. Consequently, the gene to be inserted must either be fitted with a new, host-accommodated transcription and translation regulating DNA sequence prior to insertion or it must be inserted at a site where it will come under the control of existing transcription and translation signals in the vector DNA.

It is frequently the case that the insertion of an exogenous gene into, e.g., a circular DNA plasmid vector, is performed at a site either immediately following an extant transcription and translation signal or within an existing plasmid-borne gene coding for a rather large protein which is the subject of high degrees of expression in the host. In the latter case, the host's expression of the "fusion gene" so formed results in high levels of production of a "fusion protein" including the desired protein sequence (e.g., as an intermediate segment which can be isolated by chemical cleavage of large protein). Such procedures not only insure desired regulation and high levels of expression of the exogenous gene product but also result in a degree of protection of the desired protein product from attack by proteases endogenous to the host. Further, depending on the host organism, such procedures may allow for a kind of "piggyback" transportation of the desired protein from the host cells into the cell culture medium, eliminating the need to destroy host cells for the purpose of isolating the desired product.

While the foregoing generalized descriptions of published recombinant DNA methodologies may make the processes appear to be rather straightforward, easily performed and readily verified, it is actually the case that the DNA sequence manipulations involved are quite painstakingly difficult to perform and almost invariably characterized by very low yields of desired products.

As an example, the initial "preparation" of a gene for insertion into a vector to be used in transformation of a host microorganism can be an enormously difficult process, especially where the gene to be expressed is endogenous to a higher organism such as man. One laborious procedure practiced in the art is the systematic cloning into recombinant plasmids of the total DNA genome of the "donor" cells, generating immense "libraries" of transformed cells carrying random DNA sequence fragments which must be individually tested for expression of a product of interest. According to another procedure, total mRNA is isolated from high expression donor cells (presumptively containing multiple copies of mRNA coded for the product of interest), first "copied" into single stranded cDNA with reverse transcriptase enzymes, then into double stranded form with polymerase, and cloned. The procedure again generates a library of transformed cells somewhat smaller than a total genome library which may include the desired gene copies free of non-transcribed "introns" which can significantly interfere with expression by a host microorganism. The above-noted time-consuming gene isolation procedures were in fact employed in published recombinant DNA

procedures for obtaining microorganism expression of several proteins, including rat proinsulin [Ullrich, et al., *Science*, **196**, pp. 1313-1318 (1977)], human fibroblast interferon [Goedell, et al., *Nucleic Acids Research*, **8**, pp. 4087-4094 (1980)], mouse  $\beta$ -endorphin [Shine, et al., *Nature*, **285**, pp. 456-461 (1980)] and human leukocyte interferon [Goedell, et al., *Nature*, **287**, pp. 411-416 (1980); and Goedell, et al., *Nature*, **290**, pp. 20-26 (1981)].

Whenever possible, the partial or total manufacture of genes of interest from nucleotide bases constitutes a much preferred procedure for preparation of genes to be used in recombinant DNA methods. A requirement for such manufacture is, of course, knowledge of the correct amino acid sequence of the desired polypeptide. With this information in hand, a generative DNA sequence code for the protein (i.e., a properly ordered series of base triplet codons) can be planned and a corresponding synthetic, double stranded DNA segment can be constructed. A combination of manufacturing and cDNA synthetic methodologies is reported to have been employed in the generation of a gene for human growth hormone. Specifically, a manufactured linear double stranded DNA sequence of 72 nucleotide base pairs (comprising codons specifying the first 24 amino acids of the desired 191 amino acid polypeptide) was ligated to a cDNA-derived double strand coding for amino acids Nos. 25-191 and inserted in a modified pBR322 plasmid at a locus controlled by a *lac* promoter/regulator sequence [Goedell, et al., *Nature*, **281**, pp. 544-548 (1981)].

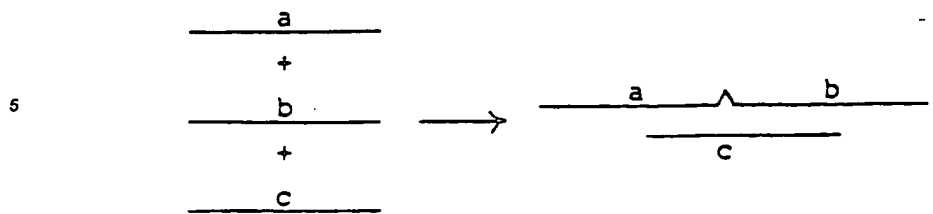
Completely synthetic procedures have been employed for the manufacture of genes coding for relatively "short" biologically functional polypeptides, such as human somatostatin (14 amino acids) and human insulin (2 polypeptide chains of 21 and 30 amino acids, respectively).

In the somatostatin gene preparative procedure [Itakura, et al., *Science*, **198**, pp. 1056-1063 (1977)] a 52 base pair gene was constructed wherein 42 base pairs represented the codons specifying the required 14 amino acids and an additional 10 base pairs were added to permit formation of "sticky-end" single stranded terminal regions employed for ligating the structural gene into a microorganism transformation vector. Specifically, the gene was inserted close to the end of a  $\beta$ -galactosidase enzyme gene and the resultant fusion gene was expressed as a fusion protein from which somatostatin was isolated by cyanogen bromide cleavage. Manufacture of the human insulin gene, as noted above, involved preparation of genes coding for a 21 amino acid chain and for a 30 amino acid chain. Eighteen deoxyoligonucleotide fragments were combined to make the gene for the longer chain, and eleven fragments were joined into a gene for the shorter chain. Each gene was employed to form a fusion gene with a  $\beta$ -galactosidase gene and the individually expressed polypeptide chains were enzymatically isolated and linked to form complete insulin molecules. [Goedell, et al., *Proc. Nat. Acad. Sci. U.S.A.*, **76**, pp. 106-110 (1979).]

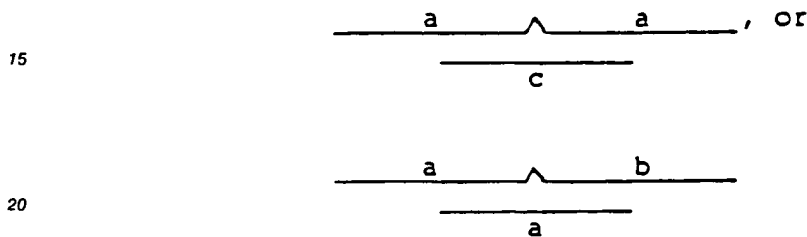
In each of the above procedures, deoxyoligonucleotide segments were prepared, and then sequentially ligated according to the following general procedure. [See, e.g., Agarwal, et al., *Nature*, **227**, pp. 1-7 (1970) and Khorana, *Science*, **203**, pp. 614-675 (1979)]. An initial "top" (i.e., 5'-3' polarity) deoxyoligonucleotide segment is enzymatically joined to a second "top" segment. Alignment of these two "top" strands is made possible using a "bottom" (i.e., 3' to 5' polarity) strand having a base sequence complementary to half of the first top strand and half of the second top strand. After joining, the uncomplemented bases of the top strands "protrude" from the duplex portion formed. A second bottom strand is added which includes the five or six base complement of a protruding top strand, plus an additional five or six bases which then protrude as a bottom single stranded portion. The two bottom strands are then joined. Such sequential additions are continued until a complete gene sequence is developed, with the total procedure being very time-consuming and highly inefficient.

The time-consuming characteristics of such methods for total gene synthesis are exemplified by reports that three months' work by at least four investigators was needed to perform the assembly of the two "short", insulin genes previously referred to. Further, while only relatively small quantities of any manufactured gene are needed for success of vector insertion, the above synthetic procedures have such poor overall yields (on the order of 20% per ligation) that the eventual isolation of even minute quantities of a selected short gene is by no means guaranteed with even the most scrupulous adherence to prescribed methods. The maximum length gene which can be synthesized is clearly limited by the efficiency with which the individual short segments can be joined. If  $n$  such ligation reactions are required and the yield of each such reaction is  $y$ , the quantity of correctly synthesized genetic material obtained will be proportional to  $y^n$ . Since this relationship is exponential in nature, even a small increase in the yield per ligation reaction will result in a substantial increase in the length of the largest gene that may be synthesized.

Inefficiencies in the above-noted methodology are due in large part to the formation of undesired intermediate products. As an example, in an initial reaction forming annealed top strands associated with a bottom, "template" strand, the desired reaction may be,



10 but the actual products obtained may be



or the like. Further, the longer the individual deoxyoligonucleotides are, the more likely it is that they will form thermodynamically stable self-associations such as "hairpins" or aggregations.

Proposals for increasing synthetic efficiency have not been forthcoming and it was recently reported that, "With the methods now available, however, it is not economically practical to synthesize genes for peptides longer than about 30 amino acid units, and many clinically important proteins are much longer". [Aharonowitz, et al. *Scientific American*, 245, No. 3, pp. 140-152, at p. 151 (1981).]

An illustration of the "economic practicalities" involved in large gene synthesis is provided by the recent publication of "successful" efforts in the total synthesis of a human leukocyte interferon gene. [Edge, et al., *Nature*, 292, pp. 756-782 (1981).] Briefly summarized, 67 different deoxyoligonucleotides containing about 15 bases were synthesized and joined in the "50 percent overlap" procedure of the type noted above to form eleven short duplexes. These, in turn were assembled into four longer duplexes which were eventually joined to provide a 514 base pair gene coding for the 166 amino acid protein. The procedure, which the authors characterize as "rapid", is reliably estimated to have consumed nearly a year's effort by five workers and the efficiency of the assembly strategy was clearly quite poor. It may be noted, for example, that while 40 pmole of each of the starting 67 deoxyoligonucleotides was prepared and employed to form the eleven intermediate-sized duplexes, by the time assembly of the four large duplexes was achieved, a yield of only about 0.01 pmole of the longer duplexes could be obtained for use in final assembly of the whole gene.

Another aspect of the practice of recombinant DNA techniques for the expression, by microorganisms, of proteins of industrial and pharmaceutical interest is the phenomenon of "codon preference". While it was earlier noted that the existing machinery for gene expression in genetically transformed host cells will "operate" to construct a given desired product, levels of expression attained in a microorganism can be subject to wide variation, depending in part on specific alternative forms of the amino acid-specifying genetic code present in an inserted exogenous gene. A "triplet" codon of four possible nucleotide bases can exist in 64 variant forms. That these forms provide the message for only 20 different amino acids (as well as transcription initiation and termination) means that some amino acids can be coded for by more than one codon. Indeed, some amino acids have as many as six "redundant", alternative codons while some others have a single, required codon. For reasons not completely understood, alternative codons are not at all uniformly present in the endogenous DNA of differing types of cells and there appears to exist a variable natural hierarchy or "preference" for certain codons in certain types of cells.

As one example, the amino acid leucine is specified by any of six DNA codons including CTA, CTC, CTG, CTT, TTA, and TTG (which correspond, respectively, to the mRNA codons, CUA, CUC, CUG, CUU, UUA and UUG). Exhaustive analysis of genome codon frequencies for microorganisms has revealed endogenous DNA of *E. coli* bacteria most commonly contains the CTG leucine-specifying codon, while the DNA of yeasts and slime molds most commonly includes a TTA leucine-specifying codon. In view of this

hierarchy, it is generally held that the likelihood of obtaining high levels of expression of a leucine-rich polypeptide by an *E. coli* host will depend to some extent on the frequency of codon use. For example, a gene rich in TTA codons will in all probability be poorly expressed in *E. coli*, whereas a CTG rich gene will probably highly express the polypeptide. In a like manner, when yeast cells are the projected transformation host cells for expression of a leucine-rich polypeptide, a preferred codon for use in an inserted DNA would be TTA. See, e.g., Grantham, et al. *Nucleic Acids Research*, 8, pp. r49-r62 (1980); Grantham, et al., *Nucleic Acids Research*, 8, pp. 1893-1912 (1980); and, Grantham, et al., *Nucleic Acids Research*, 9, pp. r43-r74 (1981).

The implications of codon preference phenomena on recombinant DNA techniques are manifest, and the phenomenon may serve to explain many prior failures to achieve high expression levels for exogenous genes in successfully transformed host organisms -- a less "preferred" codon may be repeatedly present in the inserted gene and the host cell machinery for expression may not operate as efficiently. This phenomenon directs the conclusion that wholly manufactured genes which have been designed to include a projected host cell's preferred codons provide a preferred form of foreign genetic material for practice of recombinant DNA techniques. In this context, the absence of procedures for rapid and efficient total gene manufacture which would permit codon selection is seen to constitute an even more serious roadblock to advances in the art.

Of substantial interest to the background of the present invention is the state of the art with regard to the preparation and use of a class of biologically active substances, the interferons (IFNs). Interferons are secreted proteins having fairly well-defined antiviral, antitumor and immunomodulatory characteristics. See, e.g., Gray, et al., *Nature*, 295, pp. 503-508 (1982) and Edge, et al., *supra*, and references cited therein.

On the basis of antigenicity and biological and chemical properties, human interferons have been grouped into three major classes: IFN- $\alpha$  (leukocyte), IFN- $\beta$  (fibroblast) and IFN- $\gamma$  (immune). Considerable information has accumulated on the structures and properties of the virus-induced acid-stable interferons (IFN- $\alpha$  and  $\beta$ ). These have been purified to homogeneity and at least partial amino acid sequences have been determined. Analyses of cloned cDNA and gene sequences for IFN- $\beta$ , and the IFN- $\alpha$  multigene family have permitted the deduction of the complete amino acid sequences of many of the interferons. In addition, efficient synthesis of IFN- $\beta$ , and several IFN- $\alpha$ s in *E. coli*, and IFN- $\alpha$ , in yeast, have now made possible the purification of large quantities of these proteins in biologically active form.

Much less information is available concerning the structure and properties of IFN- $\gamma$ , an interferon generally produced in cultures of lymphocytes exposed to various mitogenic stimuli. It is acid labile and does not cross-react with antisera prepared against IFN- $\alpha$  or IFN- $\beta$ . A broad range of biological activities have been attributed to IFN- $\gamma$  including potentiation of the antiviral activities of IFN- $\alpha$  and - $\beta$ , from which it differs in terms of its virus and cell specificities and the antiviral mechanisms induced. *In vitro* studies performed with crude preparations suggest that the primary function of IFN- $\gamma$  may be as an immunoregulatory agent. The antiproliferative effect of IFN- $\gamma$  on transformed cells has been reported to be 10 to 100-fold greater than that of IFN- $\alpha$  or - $\beta$ , suggesting a potential use in the treatment of neoplasia. Murine IFN- $\gamma$  preparations have been shown to have significant antitumor activity against mouse sarcomas.

It has recently been reported (Gray, et al., *supra*) that a recombinant plasmid containing a cDNA sequence coding for human IFN- $\gamma$  has been isolated and characterized. Expression of this sequence in *E. coli* and cultured monkey cells is reported to give rise to a polypeptide having the properties of authentic human IFN- $\gamma$ . In the publication, the cDNA sequence and the deduced 146 amino acid sequence of the "mature" polypeptide, exclusive of the putative leader sequence, is as follows:

1 10  
 Cys-Tyr-Cys-Gln-Asp-Pro-Tyr-Val-Lys-Glu-Ala-Glu-Asn-Leu-  
 5 TGT TAC TGC CAG CAG CAA TAT GTA AAA GAA GCA GAA AAC CTT  
 20  
 Lys-Lys-Tyr-Phe-Asn-Ala-Gly-His-Ser-Asp-Val-Ala-Asp-Asn-  
 10 AAG AAA TAT TTT AAT GCA GGT CAT TCA GAT GTA GCG GAT AAT  
 30 40  
 Gly-Thr-Leu-Phe-Leu-Gly-Ile-Leu-Lys-Asn-Trp-Lys-Glu-Glu-  
 15 GGA ACT CTT TTC TTA GGC ATT TTG AAG AAT TGG AAA GAG GAG  
 50  
 Ser-Asp-Arg-Lys-Ile-Met-Gln-Ser-Gln-Ile-Val-Ser-Phe-Tyr-  
 20 AGT GAC AGA AAA ATA ATG CAG AGC CAA ATT GTC TCC TTT TAC  
 60 70  
 Phe-Lys-Leu-Phe-Lys-Asn-Phe-Lys-Asp-Asp-Gln-Ser-Ile-Gln-  
 TTC AAA CTT TTT AAA AAC TTT AAA GAT GAC CAG AGC ATC CAA  
 25 80  
 Lys-Ser-Val-Glu-Thr-Ile-Lys-Glu-Asp-Met-Asn-Val-Lys-Phe-  
 AAG AGT GTG GAG ACC ATC AAG GAA GAC ATG AAT GTC AAG TTT  
 30 90  
 Phe-Asn-Ser-Asn-Lys-Lys-Lys-Arg-Asp-Asp-Phe-Glu-Lys-Leu-  
 TTC AAT AGC AAC AAA AAG AAA CGA GAT GAC TTC GAA AAG CTG  
 100 110  
 35 Thr-Asn-Tyr-Ser-Val-Thr-Asp-Leu-Asn-Val-Gln-Arg-Lys-Ala-  
 ACT AAT TAT TCG GTA ACT GAC TTG AAT GTC CAA CGC AAA GCA  
 120  
 40 Ile-His-Glu-Leu-Ile-Gln-Val-Met-Ala-Glu-Leu-Ser-Pro-Ala-  
 ATA CAT GAA CTC CTC ATC CAA ATG GCT GAA CTG TCG CAA GCA  
 130 140  
 45 Ala-Lys-Thr-Gly-Lys-Arg-Lys-Arg-Ser-Gln-Met-Leu-Phe-Gln-  
 GCT AAA ACA GGG AAG CGA AAA AGG AGT CAG ATG CTG TTT CAA  
 50 146  
 Gly-Arg-Arg-Ala-Ser-Gln  
 GGT CGA AGA GCA TCC CAG.

55 In a previous publication of the sequence, arginine, rather than glutamine, was specified at position 140 in the sequence. (Unless otherwise indicated, therefore, reference to "human immune interferon" or, simply "IFN- $\gamma$ " shall comprehend both the [Arg<sup>140</sup>] and [Gln<sup>140</sup>] forms.)

The above-noted wide variations in biological activities of various interferon types makes the construction of synthetic polypeptide analogs of the interferons of paramount significance to the full development of the therapeutic potential of this class of compounds. Despite the advantages in isolation of quantities of interferons which have been provided by recombinant DNA techniques to date, practitioners in this field have not been able to address the matter of preparation of synthetic polypeptide analogs of the interferons with any significant degree of success.

Put another way, the work of Gray, et al., *supra*, in the isolation of a gene coding for IFN- $\gamma$  and the extensive labors of Edge, et al., *supra*, in providing a wholly manufactured IFN- $\alpha_1$  gene provide only genetic materials for expression of single, very precisely defined, polypeptide sequences. There exist no procedures (except, possibly, for site specific mutagenesis) which would permit microbial expression of large quantities of human IFN- $\gamma$  analogs which differed from the "authentic" polypeptide in terms of the identity or location of even a single amino acid. In a like manner, preparation of an IFN- $\alpha_1$  analog which differed by one amino acid from the polypeptide prepared by Edge, et al., *supra*, would appear to require an additional year of labor in constructing a whole new gene which varied in terms of a single triplet codon. No means is readily available for the excision of a fragment of the subject gene and replacement with a fragment including the coding information for a variant polypeptide sequence. Further, modification of the reported cDNA-derived and manufactured DNA sequences to vary codon usage is not an available "option".

Indeed, the only report of the preparation of variant interferon polypeptide species by recombinant DNA techniques has been in the context of preparation and expression of "hybrids" of human genes for IFN- $\alpha_1$  and IFN- $\alpha_2$  [Weck, et al., *Nucleic Acids Research*, 9, pp. 6153-6168 (1981) and Streuli, et al., *Proc. Nat. Acad. Sci. U.S.A.*, 78, pp. 2848-2852 (1981)]. The hybrids obtained consisted of the four possible combinations of gene fragments developed upon finding that two of the eight human (cDNA-derived) genes fortuitously included only once within the sequence, base sequences corresponding to the restriction endonuclease cleavage sites for the bacterial endonucleases, PvuII and BglII.

There exists, therefore, a substantial need in the art for more efficient procedures for the total synthesis from nucleotide bases of manufactured DNA sequences coding for large polypeptides such as the interferons. There additionally exists a need for synthetic methods which will allow for the rapid construction of variant forms of synthetic sequences such as will permit the microbial expression of synthetic polypeptides which vary from naturally occurring forms in terms of the identity and/or position of one or more selected amino acids.

#### BRIEF SUMMARY

The present invention provides novel, rapid and highly efficient procedures for the total synthesis of linear, double stranded DNA sequences in excess of about 200 nucleotide base pairs in length, which sequences may comprise entire structural genes capable of directing the synthesis of a wide variety of polypeptides of interest.

According to the invention, linear, double stranded DNA sequences of a length in excess of about 200 base pairs and coding for expression of a predetermined continuous sequence of amino acids within a selected host microorganism transformed by a selected DNA vector including the sequence, are chemically synthesized by a method comprising:

(a) preparing two or more different, subunit, linear, double stranded DNA sequences of about 100 or more base pairs in length for assembly in a selected assembly vector,

each different subunit DNA sequence prepared comprising a series of nucleotide base codons coding for a different continuous portion of said predetermined sequence of amino acids to be expressed,

one terminal region of a first of said subunits comprising a portion of a base sequence which provides a recognition site for cleavage by a first restriction endonuclease, which recognition site is entirely present either once or not at all in said selected assembly vector upon insertion of the subunit therein,

one terminal region of a second of said subunits comprising a portion of a base sequence which provides a recognition site for cleavage by a second restriction endonuclease other than said first endonuclease, which recognition site is entirely present once or not at all in said selected assembly vector upon insertion of the subunit therein,

at least one-half of all remaining terminal regions of subunits comprising a portion of a recognition site (preferably a palindromic six base recognition site) for cleavage by a restriction endonuclease other than said first and second endonucleases, which recognition site is entirely present once and only once in said selected assembly vector after insertion of all subunits thereinto; and

(b) serially inserting each of said subunit DNA sequences prepared in step (a) into the selected assembly vector and effecting the biological amplification of the assembly vector subsequent to each insertion, thereby to form a DNA vector including the desired DNA sequence coding for the predetermined continuous amino acid sequence and wherein the desired DNA sequence assembled includes at least one unique, preferably palindromic six base, recognition site for restriction endonuclease cleavage at an intermediate position therein.

The above general method preferably further includes the step of isolating the desired DNA sequence from the assembly vector preferably to provide one of the class of synthesised DNA sequences having at least one unique palindromic six base recognition site for restriction endonuclease cleavage at an intermediate position therein. A sequence so isolated may then be inserted in a different, "expression" vector and direct expression of the desired polypeptide by a microorganism which is the same as or different from that in which the assembly vector is amplified. In other preferred embodiments of the method: at least three different subunit DNA sequences are prepared in step (a) and serially inserted into said selected assembly vector in step (b) and the desired synthesised DNA sequence obtained includes at least two unique palindromic six base recognition sites for restriction endonuclease cleavage at intermediate positions therein; the DNA sequence synthesized comprises an entire structural gene coding for a biologically active polypeptide; and, in the DNA sequence synthesised, the sequence of nucleotide bases includes one or more codons selected, from among alternative codons specifying the same amino acid, on the basis of preferential expression characteristics of the codon in said selected host microorganism.

Illustratively provided by the present invention are synthesised genes coding for the synthesis of human immune interferon (IFN- $\gamma$ ) and biologically functional analog polypeptides which differ from human immune interferon in terms of the identity and/or location of one or more amino acids. Also provided are synthesised genes coding for synthesis of human leukocyte interferon of the F subtype ("L $\alpha$ IFN-F" or "IFN- $\alpha$ F") and analogs thereof, along with consensus human leukocyte interferons.

DNA subunit sequences for use in practice of the methods of the invention are preferably synthesized from nucleotide bases according to the methods disclosed in co-owned, concurrently-filed U.S. Patent Application Serial No. 375,493, by Yitzhak Stabinsky (WO83/04029), entitled "Manufacture and Expression of Structural Genes" (Attorney's Docket No. 6250). Briefly summarized the general method comprises the steps of:

- (1) preparing two or more different, linear, duplex DNA strands, each duplex strand including a double stranded region of 12 or more selected complementary base pairs and further including a top single stranded terminal sequence of from 3 to 7 selected bases at one end of the strand and/or a bottom single stranded terminal sequence of from 3 to 7 selected bases at the other end of the strand, each single stranded terminal sequence of each duplex DNA strand comprising the entire base complement of at most one single stranded terminal sequence of any other duplex DNA strand prepared; and
- (2) annealing each duplex DNA strand prepared in step (1) to one or two different duplex strands prepared in step (1) having a complementary single stranded terminal sequence, thereby to form a single continuous double stranded DNA sequence which has a duplex region of at least 27 selected base pairs including at least 3 base pairs formed by complementary association of single stranded terminal sequences of duplex DNA strands prepared in step (1) and which has from 0 to 2 single stranded top or bottom terminal regions of from 3 to 7 bases.

In the preferred general process for subunit manufacture, at least three different duplex DNA strands are prepared in step (1) and all strands so prepared are annealed concurrently in a single annealing reaction mixture to form a single continuous double stranded DNA sequence which has a duplex region of at least 42 selected base pairs including at least two non-adjacent sets of 3 or more base pairs formed by complementary association of single stranded terminal sequences of duplex strands prepared in step (1).

The duplex DNA strand preparation step (1) of the preferred subunit manufacturing process preferably comprises the steps of:

- (a) constructing first and second linear deoxyoligonucleotide segments having 15 or more bases in a selected linear sequence, the linear sequence of bases of the second segment comprising the total complement of the sequence of bases of the first segment except that at least one end of the second segment shall either include an additional linear sequence of from 3 to 7 selected bases beyond those fully complementing the first segment, or shall lack a linear sequence of from 3 to 7 bases complementary to a terminal sequence of the first segment, provided, however, that the second segment shall not have an additional sequence of bases or be lacking a sequence of bases at both of its ends; and,
- (b) combining the first and second segments under conditions conducive to complementary association between segments to form a linear, duplex DNA strand.

The sequence of bases in the double stranded DNA subunit sequences formed preferably includes one or more triplet codons selected from among alternative codons specifying the same amino acid on the basis of preferential expression characteristics of the codon in a projected host microorganism, such as yeast cells or bacteria, especially E. coli bacteria.

5 Also provided by the present invention are improvements in methods and materials for enhancing levels of expression of selected exogenous genes in E. coli host cells. Briefly stated, expression vectors are constructed to include selected DNA sequences upstream of polypeptide coding regions which selected sequences are duplicative of ribosome binding site sequences extant in genomic E. coli DNA associated with highly expressed endogenous polypeptides. A presently preferred selected sequence is duplicative of  
10 the ribosome binding site sequence associated with E. coli expression of outer membrane protein F ("OMP-F").

Other aspects and advantages of the present invention will be apparent upon consideration of the following detailed description thereof.

#### 15 DETAILED DESCRIPTION

As employed herein, the term "synthesised" as applied to a DNA sequence or gene shall designate a product either totally chemically synthesized by assembly of nucleotide bases or derived from the biological replication of a product thus chemically synthesized. As such, the term is exclusive of products "syn-  
20 thesized" by cDNA methods or genomic cloning methodologies which involve starting materials which are of biological origin. Table I below sets out abbreviations employed herein to designate amino acids and includes IUPAC-recommended single letter designations.

25

30

35

40

45

50

55

TABLE I

| <u>Amino Acid</u> | <u>Abbreviation</u> | <u>IUPAC Symbol</u> |
|-------------------|---------------------|---------------------|
| Alanine           | Ala                 | A                   |
| Cysteine          | Cys                 | C                   |
| Aspartic acid     | Asp                 | D                   |
| Glutamic acid     | Glu                 | E                   |
| Phenylalanine     | Phe                 | F                   |
| Glycine           | Gly                 | G                   |
| Histidine         | His                 | H                   |
| Isoleucine        | Ile                 | I                   |
| Lysine            | Lys                 | K                   |
| Leucine           | Leu                 | L                   |
| Methionine        | Met                 | M                   |
| Asparagine        | Asn                 | N                   |
| Proline           | Pro                 | P                   |
| Glutamine         | Gln                 | Q                   |
| Arginine          | Arg                 | R                   |
| Serine            | Ser                 | S                   |
| Threonine         | Thr                 | T                   |
| Valine            | Val                 | V                   |
| Tryptophan        | Trp                 | W                   |
| Tyrosine          | Tyr                 | Y                   |

The following abbreviations shall be employed for nucleotide bases: A for adenine; G for guanine; T for thymine; U for uracil; and C for cytosine. For ease of understanding of the present invention, Table II and II below provide tabular correlations between the 64 alternate triplet nucleotide base codons of DNA and the 20 amino acids and transcription termination ("stop") functions specified thereby. In order to determine the corresponding correlations for RNA, U is substituted for T in the tables.

TABLE II

| 5  | FIRST<br>POSITION | SECOND POSITION |     |      |      | THIRD<br>POSITION |
|----|-------------------|-----------------|-----|------|------|-------------------|
|    |                   | T               | C   | A    | G    |                   |
| 10 | T                 | Phe             | Ser | Tyr  | Cys  | T                 |
|    |                   | Phe             | Ser | Tyr  | Cys  | C                 |
|    |                   | Leu             | Ser | Stop | Stop | A                 |
| 15 |                   | Leu             | Ser | Stop | Trp  | G                 |
|    | C                 | Leu             | Pro | His  | Arg  | T                 |
|    |                   | Leu             | Pro | His  | Arg  | C                 |
| 20 |                   | Leu             | Pro | Gln  | Arg  | A                 |
|    |                   | Leu             | Pro | Gln  | Arg  | G                 |
| 25 | A                 | Ile             | Thr | Asn  | Ser  | T                 |
|    |                   | Ile             | Thr | Asn  | Ser  | C                 |
|    |                   | Ile             | Thr | Lys  | Arg  | A                 |
| 30 |                   | Met             | Thr | Lys  | Arg  | G                 |
|    | G                 | Val             | Ala | Asp  | Gly  | T                 |
|    |                   | Val             | Ala | Asp  | Gly  | C                 |
| 35 |                   | Val             | Ala | Glu  | Gly  | A                 |
|    |                   | Val             | Ala | Glu  | Gly  | G                 |

40

45

50

55

TABLE III

| 5  | <u>Amino Acid</u> | <u>Specifying Codon(s)</u>   |
|----|-------------------|------------------------------|
|    | (A) Alanine       | GCT, GCC, GCA, GCG           |
| 10 | (C) Cysteine      | TGT, TGC                     |
|    | (D) Aspartic acid | GAT, GAC                     |
|    | (E) Glutamic acid | GAA, GAG                     |
| 15 | (F) Phenylalanine | TTT, TTC                     |
|    | (G) Glycine       | GGT, GGC, GGA, GGG           |
|    | (H) Histidine     | CAT, CAC                     |
|    | (I) Isoleucine    | ATT, ATC, ATA                |
| 20 | (K) Lysine        | AAA, AAG                     |
|    | (L) Leucine       | TTA, TTG, CTT, CTC, CTA, CTG |
|    | (M) Methionine    | ATG                          |
| 25 | (N) Asparagine    | AAT, AAC                     |
|    | (P) Proline       | CCT, CCC, CCA, CCG           |
|    | (Q) Glutamine     | CAA, CAG                     |
| 30 | (R) Arginine      | CGT, CGC, CGA, CGG, AGA, AGG |
|    | (S) Serine        | TCT, TCC, TCA, TCG, AGT, AGC |
|    | (T) Threonine     | ACT, ACC, ACA, ACG           |
|    | (V) Valine        | GTT, GTC, GTA, GTG           |
| 35 | (W) Tryptophan    | TGG                          |
|    | (Y) Tyrosine      | TAC, TAT                     |
|    | STOP              | TAA, TAG, TGA                |

A "palindromic" recognition site for restriction endonuclease cleavage of double stranded DNA is one which displays "left-to-right and right-to-left" symmetry between top and bottom base complements, i.e., where "readings" of complementary base sequences of the recognition site from 5' to 3' ends are identical.

45 Examples of palindromic six base recognition sites for restriction endonuclease cleavage include the sites for cleavage by HindIII wherein top and bottom strands read from 5' to 3' as AAGCTT. A non-palindromic six base restriction site is exemplified by the site for cleavage by EcoP15, the top strand of which reportedly reads CAGCAG. The bottom strand base complement, when read 5' to 3' is CTGCTG. Essentially by definition, restriction sites comprising odd numbers of bases (e.g., 5, 7) are non-palindromic.

50 Certain endonucleases will cleave at variant forms of a site, which may be palindromic or not. For example, XhoI will recognize a site which reads (any purine)GATC(any pyrimidine) including the palindromic sequence AGATCT and the non-palindromic sequence GGATCT. Referring to the previously-noted "BRL Restriction Endonuclease Reference Chart," endonucleases recognizing six base palindromic sites exclusively include BbrI, CuiI, Hin173, Ein91R, HinBIII, HinBIII, HindIII, HinfII, HsuI, BglII, StuI, RruI, ClaI, AvalI,

55 PvuII, SmaI, XmaI, EcoI, SacII, SboI, SbrI, ShyI, SstII, TglI, AvrII, PvuI, RshI, RspI, XniI, XorII, XmaIII, BlnI, MsiI, SclI, SexI, SgoI, SliI, SliI, SpaI, XhoI, XpaI, Bce170, Bsu1247, PstI, SalPI, XmaII, XorI, EcoRI, Rsh630I, SacI, SstI, SphI, BamHI, BamKI, BamNI, BamFI, BstI, KpnI, Sall, XmaI, HpaI, XbaI, AtuCI, BclI, CpeI, SstIV, Aosl, MstI, Ball, AsuII, and MlaI. Endonucleases which recognize only non-palindromic six base

sequences exclusively include TthIII, EcoP15, Aval, and Avrl. Endonucleases recognizing both palindromic and non-palindromic six base sequences include HaeI, HgiAI, AclI, AoslI, AsuII, AclI, ChulI, HincII, HindII, MnlI, XhoII, HaeII, HincII, NgoI, and EcoRI'.

Upon determination of the structure of a desired polypeptide to be produced, practice of the present invention involves: preparation of two or more different specific, continuous double stranded DNA subunit sequences of 100 or more base pairs in length and having terminal portions of the proper configuration; serial insertion of subunits into a selected assembly vector with intermediate amplification of the hybrid vectors in a selected host organism; use of the assembly vector (or an alternate, selected "expression" vector including the DNA sequence which has been manufactured from the subunits) to transform a suitable, selected host; and, isolating polypeptide sequences expressed in the host organism. In its most efficient forms, practice of the invention involves using the same vector for assembly of the manufactured sequence and for large scale expression of the polypeptide. Similarly, the host microorganism employed for expression will ordinarily be the same as employed for amplifications performed during the subunit assembly process.

The manufactured DNA sequence may be provided with a promoter/regulator region for autonomous control of expression or may be incorporated into a vector in a manner providing for control of expression by a promoter/regulator sequence extant in the vector. Manufactured DNA sequences of the invention may suitably be incorporated into existing plasmid-borne genes (e.g.,  $\beta$ -galactosidase) to form fusion genes coding for fusion polypeptide products including the desired amino acid sequences coded for by the manufactured DNA sequences.

In practice of the invention in its preferred forms, polypeptides produced may vary in size from about 65 or 70 amino acids up to about 200 or more amino acids. High levels of expression of the desired polypeptide by selected transformed host organisms is facilitated through the manufacture of DNA sequences which include one or more alternative codons which are preferentially expressed by the host.

Manufacture of double stranded subunit DNA sequences of 100 to 200 base pairs in length may proceed according to prior art assembly methods previously referred to, but is preferably accomplished by means of the rapid and efficient procedures disclosed in the aforementioned U.S. Application S.N. 375,493 by Stabinsky (WO 83/04029) and used in certain of the following examples of actual practice of the present invention. Briefly put, these procedures involve the assembly from deoxyoligonucleotides of two or more different, linear, duplex DNA strands each including a relatively long double stranded region along with a relatively short single stranded region on one or both opposing ends of the double strand. The double stranded regions are designed to include codons needed to specify assembly of an initial, or terminal or intermediate portion of the total amino acid sequence of the desired polypeptide. Where possible, alternative codons preferentially expressed by a projected host (e.g., *E. coli*) are employed. Depending on the relative position to be assumed in the finally assembled subunit DNA sequence, the single stranded region(s) of the duplex strands will include a sequence of bases which, when complemented by bases of other duplex strands, also provide codons specifying amino acids within the desired polypeptide sequence.

Duplex strands formed according to this procedure are then enzymatically annealed to the one or two different duplex strands having complementary short, single stranded regions to form a desired continuous double stranded subunit DNA sequence which codes for the desired polypeptide fragment.

High efficiencies and rapidity in total sequence assembly are augmented in such procedures by performing a single annealing reaction involving three or more duplex strands, the short, single stranded regions of which constitute the base complement of at most one other single stranded region of any other duplex strand. Providing all duplex strands formed with short single stranded regions which uniquely complement only one of the single stranded regions of any other duplex is accomplished by alternative codon selection within the context of genetic code redundancy, and preferably also in the context of codon preferences of the projected host organism.

The following description of the manufacture of a hypothetical long DNA sequence coding for a hypothetical polypeptide will serve to graphically illustrate practice of the invention, especially in the context of formation of proper terminal sequences on subunit DNA sequences.

A biologically active polypeptide of interest is isolated and its amino acids are sequenced to reveal a constitution of 100 amino acid residues in a given continuous sequence. Formation of a manufactured gene for microbial expression of the polypeptide will thus require assembly of at least 300 base pairs for insertion into a selected viral or circular plasmid DNA vector to be used for transformation of a selected host organism.

A preliminary consideration in construction of the manufactured gene is the identity of the projected microbial host, because foreknowledge of the host allows for codon selection in the context of codon preferences of the host species. For purposes of this discussion, the selection of an *E. coli* bacterial host is

posited.

A second consideration in construction of the manufactured gene is the identity of the projected DNA vector employed in the assembly process. Selection of a suitable vector is based on existing knowledge of sites for cleavage of the vector by restriction endonuclease enzymes. More particularly, the assembly vector is selected on the basis of including DNA sequences providing endonuclease cleavage sites which will permit easy insertion of the subunits. In this regard, the assembly vector selected preferably has at least two restriction sites which occur only once (i.e., are "unique") in the vector prior to performance of any subunit insertion processes. For the purposes of this description, the selection of a hypothetical circular DNA plasmid pBR 3000 having a single EcoRI restriction site, i.e.,

-GAATTC-,  
-CTTAAG-

and a single PvuII restriction site, i.e.,

-CAGCTG-,  
-STCGAC-

is posited.

The amino acid sequence of the desired polypeptide is then analyzed in the context of determining availability of alternate codons for given amino acids (preferably in the context of codon preferences of the projected *E. coli* host). With this information in hand, two subunit DNA sequences are designed, preferably having a length on the order of about 150 base pairs -- each coding for approximately one-half of the total amino acid sequences of the desired polypeptide. For purposes of this description, the two subunits manufactured will be referred to as "A" and "B".

The methods of the present invention as applied to two such subunits, generally call for: insertion of one of the subunits into the assembly vector; amplification of the hybrid vector formed; and insertion of the second subunit to form a second hybrid including the assembled subunits in the proper sequence. Because the method involves joining the two subunits together in a manner permitting the joined ends to provide a continuous preselected sequence of bases coding for a continuous preselected sequence of amino acids, there exist certain requirements concerning the identity and sequence of the bases which make up the terminal regions of the manufactured subunits which will be joined to another subunit. Because the method calls for joining subunits to the assembly vector, there exist other requirements concerning the identity and sequence of the bases which make up those terminal regions of the manufactured subunits which will be joined to the assembly vector. Because the subunits are serially, rather than concurrently, inserted into the assembly vector (and because the methods are most beneficially practiced when the subunits can be selectively excised from assembled form to allow for alterations in selected base sequences therein), still further requirements exist concerning the identity of the bases in terminal regions of subunits manufactured. For ease of understanding in the following discussion of terminal region characteristics, the opposing terminal regions of subunits A and B are respectively referred to as A-1 and A-2, and B-1 and B-2, viz:

|            |            |            |            |
|------------|------------|------------|------------|
| <u>B-2</u> | <u>B-1</u> | <u>A-2</u> | <u>A-1</u> |
| <u>B</u>   |            | <u>A</u>   |            |

Assume that an assembly strategy is developed wherein subunit A is to be inserted into pBR3000 first, with terminal region A-1 to be ligated to the vector at the EcoRI restriction site. In the simplest case, the terminal region is simply provided with an EcoRI "sticky end", i.e., a single strand of four bases (-AATT- or -TTAA-) which will complement a single stranded sequence formed upon EcoRI digestion of pBR3000. This will allow ligation of terminal region A-1 to the vector upon treatment with ligase enzyme. Unless the single strand at the end of terminal region A-1 is preceded by an appropriate base pair

(e.g., 5'-G-  
3'-CTTAA-),

the entire recognition site will not be reconstituted upon ligation to the vector. Whether or not the EcoRI recognition site is reconstituted upon ligation (i.e., whether or not there will be 0 or 1 EcoRI sites remaining after insertion of subunit A into the vector) is at the option of the designer of the strategy. Alternatively, one may construct the terminal region A-1 of subunit A to include a complete set of base pairs providing a recognition site for some other endonuclease, hypothetically designated "XXX", and then add on portions of the EcoRI recognition site as above to provide an EcoRI "linker". To be of practical use in excising subunit A from an assembled sequence, the "XXX" site should not appear elsewhere in the hybrid plasmid formed upon insertion. The requirement for construction of terminal region A-1 is, therefore, that it comprise a portion (i.e., all or part) of a base sequence which provides a recognition site for cleavage by a restriction endonuclease, which recognition site is entirely present either once or not at all in the assembly vector upon insertion of the subunit.

Assume that terminal region B-2 of subunit B is also to be joined to the assembly vector (e.g., at the single recognition site for PvuII cleavage present on pBR3000). The requirements for construction of terminal region B-2 are the same as for construction of A-1, except that the second endonuclease enzyme in reference to which the construction of B-2 is made must be different from that with respect to which the construction of A-1 is made. If recognition sites are the same, one will not be able to separately excise segments A and B from the fully assembled sequence.

The above assumptions require, then, that terminal region A-2 is to be ligated to terminal region B-1 in the final pBR3000 hybrid. Either the terminal region A-2 or the terminal region B-1 is constructed to comprise a portion of a (preferably palindromic six base) recognition site for restriction endonuclease cleavage by hypothetical third endonuclease "YYY" which recognition site will be entirely present once and only once in the expression vector upon insertion of all subunits therein, i.e., at an intermediate position in the assemblage of subunits. There exist a number of strategies for obtaining this result. In one alternative strategy, the entire recognition site of "YYY" is contained in terminal region A-2 and the region additionally includes the one or more portions of other recognition sites for endonuclease cleavage needed to (1) complete the insertion of subunit A into the assembly vector for amplification purposes, and (2) allow for subsequent joining of subunit A to subunit B. In this case, terminal region B-1 would have at its end only the bases necessary to link it to terminal region A-2. In another alternative, the entire "YYY" recognition site is included in terminal region B-1 and B-1 further includes at its end a portion of a recognition site for endonuclease cleavage which is useful for joining subunit A to subunit B.

As another alternative, terminal region B-1 may contain at its end a portion of the "YYY" recognition site. Terminal region A-2 would then contain the entire "YYY" recognition site plus, at its end, a suitable "linker" for joining A-2 to the assembly vector prior to amplification of subunit A (e.g., a PvuII "sticky end"). After amplification of the hybrid containing subunit A, the hybrid would be cleaved with "YYY" (leaving a sticky-ended portion of the "YYY" recognition site exposed on the end of A-2) and subunit B could be inserted with its B-1 terminal region joined with the end of terminal region A-2 to reconstitute the entire "YYY" recognition site. The requirement for construction of the terminal regions of all segments (other than A-1 and B-2) is that one or the other or both (i.e., "at least half") comprise a portion (i.e., include all or part) of a recognition site for third restriction endonuclease cleavage, which recognition site is entirely present once and only once (i.e., is "unique") in said assembly vector after insertion of all subunits therein. To generate a member of the class of novel DNA sequences of the invention, the recognition site of the third endonuclease should be a six base palindromic recognition site.

While a subunit "terminal region" as referred to above could be considered to extend from the subunit end fully halfway along the subunit to its center, as a practical matter the constructions noted would ordinarily be performed in the final 10 or 20 bases. Similarly, while the unique "intermediate" recognition site in the two subunit assemblage may be up to three times closer to one end of the manufactured sequence than it is to the other, it will ordinarily be located near the center of the sequence. If, in the above description, a synthetic plan was generated calling for preparation of three subunits to be joined, the manufactured gene would include two unique restriction enzyme cleavage sites in intermediate positions at least one of which will have a palindromic six base recognition site in the class of new DNA sequences of the invention.

The significant advantages of the above-described process are manifest. Because the manufactured gene now includes one or more unique restriction endonuclease cleavage sites at intermediate positions along its length, modifications in the codon sequence of the two subunits joined at the cleavage site may be effected with great facility and without the need to re-synthesize the entire manufactured gene.

Following are illustrative examples of the actual practice of the invention in formation of manufactured genes capable of directing the synthesis of: human immune interferon (IFN- $\gamma$ ) and analogs thereof; human leukocyte interferon of the F subtype (INF- $\alpha$ F) and analogs thereof; and, multiple consensus leukocyte

interferons which, due to homology to IFN- $\alpha$ F can be named as IFN- $\alpha$ F analogs. It will be apparent from these examples that the gene manufacturing methodology of the present invention provides an overall synthetic strategy for the truly rapid, efficient synthesis and expression of genes of a length in excess of 200 base pairs within a highly flexible framework allowing for variations in the structures of products to be expressed which has not heretofore been available to investigators practicing recombinant DNA techniques.

#### EXAMPLE 1

In the procedure for construction of synthetic genes for expression of human IFN- $\gamma$  a first selection made was the choice of *E. coli* as a microbial host for eventual expression of the desired polypeptides. Thereafter, codon selection procedures were carried out in the context of *E. coli* codon preferences enumerated in the Grantham publications, supra. A second selection made was the choice of pBR322 as an expression vector and, significantly, as the assembly vector to be employed in amplification of subunit sequences. In regard to the latter factor, the plasmid was selected with the knowledge that it included single BamHI, HindIII, and Sall restriction sites. With these restriction sites and the known sequence of amino acids in human immune interferon in mind, a general plan for formation of three "major" subunit DNA sequences (IF-3, IF-2 and IF-1) and one "minor" subunit DNA sequence (IF-4) was evolved. This plan is illustrated by Table IV below.

**TABLE IV**

**IF-4**

**lipa I**

HpaI  
 A ACT AGT ACG CAA GTT CAC GTA AAA AGG GTA TCG ACA ATG TGT TAC TGC CAG  
 T TGA TCA TGC GTT CAA GTG CAT TTT TCC CAT AGC TGT TAC ACA ATG ACG GTC CTAG  
 BamHI  
 -1 1 2 3 4  
 Met Cys Tyr Cys Gln

TABLE IV (cont'd.)

IF-3

|    |     |     |     |     |     |     |     |     |     |     |     |     |     |     |     |     |
|----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|
| 4  | Gln | Asp | Pro | Tyr | Val | Lys | Glu | Ala | Glu | Asn | Leu | Lys | Lys | Tyr | Phe | Asn |
| 10 |     |     |     |     |     |     |     |     |     |     |     |     |     |     |     |     |
| 15 | A   | ATT | CGA | CAG | GAT | CCG | TAC | GTT | AAG | GAA | GCA | AAC | CTG | AAA | AAA | TAC |
| 20 |     |     |     |     |     |     |     |     |     |     |     |     |     |     |     |     |
| 25 |     |     |     |     |     |     |     |     |     |     |     |     |     |     |     |     |
| 30 |     |     |     |     |     |     |     |     |     |     |     |     |     |     |     |     |
| 35 |     |     |     |     |     |     |     |     |     |     |     |     |     |     |     |     |
| 40 |     |     |     |     |     |     |     |     |     |     |     |     |     |     |     |     |
| 45 |     |     |     |     |     |     |     |     |     |     |     |     |     |     |     |     |
| 50 |     |     |     |     |     |     |     |     |     |     |     |     |     |     |     |     |
| 55 |     |     |     |     |     |     |     |     |     |     |     |     |     |     |     |     |

Bam HI

|    |     |     |     |     |     |     |     |     |     |     |     |     |     |     |     |     |     |     |     |     |
|----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|
| 20 | Ala | Gly | His | Ser | Asp | Val | Ala | Asp | Asn | Gly | Thr | Leu | Phe | Leu | Gly | Ile | Leu | Lys | Asn | Trp |
| 30 |     |     |     |     |     |     |     |     |     |     |     |     |     |     |     |     |     |     |     |     |
| 35 |     |     |     |     |     |     |     |     |     |     |     |     |     |     |     |     |     |     |     |     |
| 40 |     |     |     |     |     |     |     |     |     |     |     |     |     |     |     |     |     |     |     |     |
| 45 |     |     |     |     |     |     |     |     |     |     |     |     |     |     |     |     |     |     |     |     |
| 50 |     |     |     |     |     |     |     |     |     |     |     |     |     |     |     |     |     |     |     |     |
| 55 |     |     |     |     |     |     |     |     |     |     |     |     |     |     |     |     |     |     |     |     |

IF-2

[illegible][illegible]

Lys Arg Asp Asp Phe Glu Lys Lys Leu Thr Asn Tyr Ser Val Thr  
100

AAA CGC GAT GAC TTC GAG AAA CTG ACT AAC TAC TCT GTT ACA  
TTT GCG CTA CTG AAG CTC TTT GAC TGA TTG ATG AGA CAA TGT CTA G

IF-1

[illegible]

Alternative form of subunit IFN-1 for use in synthesis of a manufactured gene for [Arg<sup>140</sup>]IFN-γ included the codon 5'-CGT in place of 5'-CAG (for [Gln<sup>140</sup>]) at the codon site specifying the 140th amino acid.

The codon sequence plan for the top strand of the polypeptide-specifying portion total DNA sequence synthesized was as follows:

5        5' -TGT-TAC-TGC-CAG-GAT-CCG-TAC-GTT-AAG-GAA-GCA-GAA-  
          AAC-CTG-AAA-AAA-TAC-TTC-AAC-GCA-GGC-CAC-TCC-GAC-GTA-  
          GCT-GAT-AAC-GGC-ACC-CTG-TTC-CTG-GGT-ATC-CTA-AAA-AAC-  
 10       TGG-AAA-GAG-GAA-TCC-GAC-CTG-AAG-ATC-ATG-CAG-TCT-CAA-  
          ATT-GTA-AGC-TTC-TAC-TTC-AAA-CTG-TTC-AAG-AAC-TTC-AAA-  
          GAC-GAT-CAA-TCC-ATC-CAG-AAG-AGC-GTA-GAA-ACT-ATT-AAG-  
          GAG-GAC-ATG-AAC-GTA-AAA-TCC-TTT-AAC-AGC-AAC-AAG-AAG-  
 15       AAA-CGC-GAT-GAC-TTC-GAG-AAA-CTG-ACT-AAC-TAC-TCT-GTT-  
          ACA-GAT-CTG-AAC-GTG-CAG-CGT-AAA-GCT-ATT-CAC-GAA-CTG-  
          ATC-CAA-GTT-ATG-GCT-GAA-CTG-TCT-CCT-GCG-GCA-AAG-ACT-  
 20       GGC-AAA-CGC-AAG-CGT-AGC-CAG-ATG-CTG-TTT-CAG- [or CGT] -  
          CGT-CGC-CGT-GCT-TCT-CAG .

25       In the above sequence, the control sequence bases and the initial methionine-specifying codon is not illustrated, nor are termination sequences or sequences providing a terminal Sall restriction site. Vertical lines separate top strand portions attributable to each of the subunit sequences.

The following example illustrates a preferred general procedure for preparation of deoxyoligonucleotides for use in the manufacture of DNA sequences of the invention.

#### 30       EXAMPLE 2

Oligonucleotide fragments were synthesized using a four-step procedure and several intermediate washes. Polymer bound dimethoxytrityl protected nucleoside in a sintered glass funnel was first stripped of  
 35       its 5'-protecting group (dimethoxytrityl) using 3% trichloroacetic acid in dichloromethane for 1-1/2 minutes. The polymer was then washed with methanol, tetrahydrofuran and acetonitrile. The washed polymer was then rinsed with dry acetonitrile, placed under argon and then treated in the condensation step as follows. 0.5 ml of a solution of 10 mg tetrazole in acetonitrile was added to the reaction vessel containing polymer. Then 0.5 ml of 30 mg protected nucleoside phosphoramidite in acetonitrile was added. This reaction was  
 40       agitated and allowed to react for 2 minutes. The reactants were then removed by suction and the polymer rinsed with acetonitrile. This was followed by the oxidation step wherein 1 ml of a solution containing 0.1 molar I<sub>2</sub> in 2-6-lutidine/H<sub>2</sub>O/THF, 1:2:2, was reacted with the polymer bound oligonucleotide chain for 2 minutes. Following a THF rinse capping was done using a solution of dimethylaminopyridine (6.5 g in 100 ml THF) and acetic anhydride in the proportion 4:1 for 2 minutes. This was followed by a methanol rinse  
 45       and a THF rinse. Then the cycle began again with a trichloroacetic acid in CH<sub>2</sub>Cl<sub>2</sub> treatment. The cycle was repeated until the desired oligonucleotide sequence was obtained.

The final oligonucleotide chain was treated with thiophenol dioxane, triethylamine 1:2:2, for 45 minutes at room temperature. Then, after rinsing with dioxane, methanol and diethylether, the oligonucleotide was cleaved from the polymer with concentrated ammonium hydroxide at room temperature. After decanting the  
 50       solution from the polymer, the concentrated ammonium hydroxide solution was heated at 60 °C for 16 hours in a sealed tube.

Each oligonucleotide solution was then extracted four times with 1-butanol. The solution was loaded into a 20% polyacrylamide 7 molar urea electrophoresis gel and, after running, the appropriate product DNA band was isolated.

55       Subunits were then assembled from deoxyoligonucleotides according to the general procedure for assembly of subunit IF-1.

Following the isolation of the desired 14 DNA segments, subunit IF-1 was constructed in the following manner:

1. One nanomole of each of the DNA fragments excluding segment 13 and segment 2 which contain 5' cohesive ends, were subjected to 5'-phosphorylation;

2. The complementary strands of DNA, segments 13 and 14, 11 and 12, 9 and 10, 7 and 8, 5 and 6, 3 and 4 and 1 and 2 were combined together, warmed to 90° and slowly cooled to 25°;

5 3. The resulting annealed pairs of DNA were combined sequentially and warmed to 37° and slowly cooled to 25°;

4. The concentration of ATP and DTT in the final tube containing segments 1 thru 14 was adjusted to 150  $\mu$ M and 18 mM respectively. Twenty units of T-4 DNA ligase was added to this solution and the reaction was incubated at 4° for 18 hrs;

10 5. The resulting crude product was heated to 90° for 2 min. and subjected to gel filtration on Sephadex G50/40 using 10 mM triethyl ammonium bicarbonate as the eluent;

6. The desired product was purified, following 5' phosphorylation, using an 8% polyacrylamide-TBE gel. Subunits IF-2, IF-3 and IF-4 were constructed in a similar manner.

The following example relates to: assembly of the complete human immune interferon gene from subunits IF-1, IF-2, IF-3, and IF-4; procedures for the growing, under appropriate nutrient conditions, of transformed *E. coli* cells, the isolation of human immune interferon from the cells, and the testing of biological activity of interferon so isolated.

### EXAMPLE 3

20 The major steps in the general procedure for assembly of the complete human IFN- $\gamma$  specifying genes from subunits IF-1, IF-2, and IF-3 are illustrated in Figure 1.

The 136 base hair subunit IF-1 was electroeluted from the gel, ethanol precipitated and resuspended in water at a concentration of 0.05 pmol/ $\mu$ l. Plasmid pBR322 (2.0 pmol) was digested with EcoRI and Sall, treated with phosphatase, phenol extracted, ethanol precipitated, and resuspended in water at a concentration of 0.1 pmol/ $\mu$ l. Ligation was carried out with 0.1 pmol of the plasmid and 0.2 pmol of subunit IF-1, using T-4 DNA ligase to form hybrid plasmid pINT1. *E. coli* were transformed and multiple copies of pINT1 were isolated therefrom.

The above procedure was repeated for purposes of inserting the 153 base pair subunit IF-2 to form pINF2 except that the plasmid was digested with EcoRI and BglII. The 153 base pair IF-3 subunit was similarly inserted into pINT2 during manufacture of pINT3 except that EcoRI and Hind III were used to digest the plasmid.

An IF-4 subunit was employed in the construction of the final expression vector as follows: Plasmid PVvl was purchased from Stanford University, Palo Alto, California, and digested with PvuII. Using standard procedures, an EcoRI recognition site was inserted in the plasmid at a PvuII site. Copies of this hybrid were then digested with EcoRI and HpaI to provide a 245 base pair sequence including a portion of the trp promoter/operator region. By standard procedures, IF-4 was added to the HpaI site in order to incorporate the remaining 37 base pairs of the complete trp translational initiation signal and bases providing codons for the initial four amino acids of immune interferon (Cys-Tyr-Cys-Gln). The resulting assembly was then inserted into pINT3 which had been digested with EcoRI and BamHI to yield a plasmid designated pINT- $\gamma$ -trpI7.

*E. coli* cells containing pINT- $\gamma$ -trpI7 were grown on K media in the absence of tryptophan to an O.D.<sub>600</sub> of 1. Indoleacrylic acid was added at a concentration of 20  $\mu$ g per ml and the cells were cultured for an additional 2 hours at 37°C. Cells were harvested by centrifugation and the cell pellet was resuspended in fetal calf serum buffered with HEPES (pH 8.0). Cells were lysed by one passage through a French press at 10,000 psi. The cell lysate was cleared of debris by centrifugation and the supernatant was assayed for antiviral activity by the CPE assay ["The Interferon System" Stewart, ed., Springer-Verlag, N.Y., N.Y. (1981)]. The isolated product of expression was designated  $\gamma$ -1.

This example relates to a modification in the DNA sequence of plasmid pINT- $\gamma$ -trpI7 which facilitated the use of the vector in the trp promoter-controlled expression of structural genes coding for, e.g., analogs of IFN- $\gamma$  and IFN- $\alpha$ F.

### EXAMPLE 4

55 Segment IF-4, as previously noted, had been constructed to include bases coding for an initial methionine and the first four amino acids of IFN- $\gamma$  as well as 37 base pairs (commencing at its 5' end with a HpaI blunt end) which completed at the 3' end of a trp promoter/operator sequence, including a Shine Delgarno ribosome binding sequence. It was clear that manipulations involving sequences coding IFN- $\gamma$

analogues and for polypeptides other than IFN- $\gamma$  would be facilitated if a restriction site 3' to the entire trp promoter/operator region could be established. By way of illustration, sequences corresponding to IF-4 for other genes could then be constructed without having to reconstruct the entire 37 base pairs needed to reconstitute the trp promoter/operator and would only require bases at the 5' end such as would facilitate insertion in the proper reading frame with the complete promoter/operator.

Consistent with this goal, sequence IF-4 was reconstructed to incorporate an XbaI restriction site 3' to the base pairs completing the trp promoter/operator. The construction is shown in Table V below.

TABLE V

HpaI

```

      a                                     c
AA CTA GTA CGC AAG TTC ACG TAA AAA GGG
TT GAT CAT GCG TTC AAG TGC ATT TTT CCC
      b

```

```

      XbaI      -1  1  2  3  4      BamHI
                  Met Cys Tyr Cys Gln

```

```

      e
TAT CTA GAA ATG TGT TAC TGC CAG
ATA GAT CTT TAC ACA ATG ACG GTC CTAG
      d                                     f

```

This variant form of segment IF-4 was inserted in pINT- $\gamma$ -trp17 (digested with HpaI and BamHI) to generate plasmid pINT- $\gamma$ -TXb4 from which the IFN- $\gamma$ -specifying gene could be deleted by digestion with XbaI and Sall and the entire trp promoter/operator would remain on the large fragment.

The following example relates to construction of structural analogs of IFN- $\gamma$  whose polypeptide structure differs from that of IFN- $\gamma$  in terms of the identity of location of one or more amino acids.

#### EXAMPLE 5

A first class of analogs of IFN- $\gamma$  was formed which included a lysine residue at position 81 in place of asparagine. The single base sequence change needed to generate this analog was in subunit IF-2 of Table IV in segments 35 and 36. The asparagine-specifying codon, AAC, was replaced by the lysine-specifying codon, AAG. The isolated product of expression of such a modified DNA sequence [Lys<sup>81</sup>]IFN- $\gamma$ , was designated  $\gamma$ -10.

Another class of IFN- $\gamma$  analogs consists of polypeptides wherein one or more potential glycosylation sites present in the amino acid sequence are deleted. More particularly, these consist of [Arg<sup>140</sup>]IFN- $\gamma$  or [Gln<sup>140</sup>]IFN- $\gamma$  wherein the polypeptide sequence fails to include one or more naturally occurring sequences, [(Asn or Gln)-(ANY)-(Ser or Thr)], which are known to provide sites for glycosylation of the polypeptide. One such sequence in IFN- $\gamma$  spans positions 28 through 30, (Asn-Gly-Thr), another spans positions 101 through 103 (Asn-Tyr-Ser). Preparation of an analog according to the invention with a modification at positions 28-30 involved cleavage of plasmid containing all four IFN- $\gamma$  subunits with BamHI and HindIII to delete subunit IF-3, followed by insertion of a variant of subunit IF-3 wherein the AAC codon for asparagine therein is replaced by the codon for glutamine, CAG. (Such replacement is effected by modification of deoxyoligonucleotide segment 37 to include CAG rather than AAC and of segment 38 to include GTC rather than TTG. See Table IV.) The isolated product of expression of such a modified DNA sequence, [Gln<sup>28</sup>]IFN- $\gamma$ , was designated  $\gamma$ -12. Polypeptide analogs of this type would likely not be glycosylated if expressed in yeast cells. Polypeptide analogs as so produced are not expected to differ appreciably from naturally-occurring IFN- $\gamma$  in terms of reactivity with antibodies to the natural form, or in duration of antiproliferative or

immunomodulatory pharmacological effects, but may display enhanced potency of pharmacological activity in one or more manner.

Other classes of IFN- $\gamma$  analogs consists of polypeptides wherein the [Trp<sup>39</sup>] residue is replaced by [Phe<sup>39</sup>], and/or wherein one or more of the methionine residues at amino acid positions 48, 80, 120 and 137 are replaced by, e.g., leucine, and/or wherein cysteines at amino acid positions 1 and 3 are replaced by, e.g., serine or are completely eliminated. These last-mentioned analogs may be more easily isolated upon microbial expression because they lack the capacity for formation of intermolecular disulfide bridge formation.

Replacement of tryptophan with phenylalanine at position 39 required substitution for a TGG codon in subunit IF-3 with TTC (although TTT could also have been used), effected by modification of the deoxyoligonucleotide segment 33 (TGG to TTC) and overlapping segment 36 (TGA to TAC) used to manufacture IF-3. [Phe<sup>39</sup>, Lys<sup>81</sup>]IFN- $\gamma$ , the isolated product of expression of such a modified DNA sequence (which also included the above-noted replacement of asparagine by lysine at position 81) was designated  $\gamma$ -5.

In a like manner, replacement of one or more methionines at positions 48, 80, 120, and 137, respectively, involves alteration of subunit IF-3 (with reconstruction of deoxyoligonucleotides 31, 32 and 34), subunit IF-2 (with reconstruction of deoxyoligonucleotide segments 21 and 22); and subunit IF-1 (with reconstruction of deoxyoligonucleotide segments 7 and 10 and/or 3 and 4). An analog of IFN- $\gamma$  wherein threonine replaced methionine at position 48 was obtained by modification of segment 31 in subunit IF-3 to delete the methionine-specifying codon ATG and replace it with an ACT codon. Alterations in segments 34 (TAC to TGA) were also needed to effect this change. [Thr<sup>48</sup>, Lys<sup>81</sup>]IFN- $\gamma$ , the isolated product of expression of such a modified DNA sequence (also including a lysine-specifying codon at position 81) was designated  $\gamma$ -6.

Replacement or deletions of cysteines at positions 1 and 3 involves only alteration of subunit IF-4. As a first example, modifications in construction of subunit IF-4 to replace both of the cysteine-specifying codons at positions 1 and 3 (TGT and TGC, respectively) with the serine-specifying codon, TCT, required reconstruction of only 2 segments (see e and f of Table IV). [Ser<sup>1</sup>, Ser<sup>3</sup>, Lys<sup>81</sup>]IFN- $\gamma$ , the isolated product of expression of the thus modified [Lys<sup>81</sup>]IFN- $\gamma$  DNA sequence, was designated  $\gamma$ -2. As another example, [Lys<sup>1</sup>, Lys<sup>2</sup>, Gln<sup>3</sup>, Lys<sup>81</sup>]IFN- $\gamma$ , designated  $\gamma$ -3, was obtained as an expression product of a modified construction of subunit IF-4 wherein codons AAA, AAA, and CAA respectively replaced TTG, TAC and TGC. Finally, [des-Cys<sup>1</sup>, des-Tyr<sup>2</sup>, des-Cys<sup>3</sup>, Lys<sup>81</sup>]IFN- $\gamma$ , designated  $\gamma$ -4, was obtained by means of modification of subunit IF-4 sections to

5'-ATG CAG-3'  
3'-TAC GTC-5'

in the amino acid specifying region. It should be noted that the above modifications in the initial amino acid coding regions of the gene were greatly facilitated by the construction of pINT- $\gamma$ -TXb4 in Example 4 which meant that only short sequences with XbaI and BamHI sticky ends needed to be constructed to complete the amino terminal protein coding sequence and link the gene to the complete trp promoter.

Among other classes of IFN- $\gamma$  analog polypeptide provided by the present invention are those including polypeptides which differ from IFN- $\gamma$  in terms of amino acids traditionally held to be involved in secondary and tertiary configuration of polypeptides. As an example, provision of a cysteine residue at an intermediate position in the IFN- $\gamma$  polypeptide may generate a species of polypeptide structurally facilitative of formation of intramolecular disulfide bridges between amino terminal and intermediate cysteine residues such as found in IFN- $\alpha$ . Further, insertion or deletion of prolines in polypeptides according to the invention may alter linear and bending configurations with corresponding effects on biological activity. [Lys<sup>81</sup>, Cys<sup>95</sup>]IFN- $\gamma$ , designated  $\gamma$ -9, was isolated upon expression of a DNA sequence fashioned with

5'-TCG-3'  
3'-AGC-5'

replacing

5'-TTC-3'  
3'-AAG-5'

5

in sections 17 and 18 of subunit IF-2. A DNA sequence specifying [Cys<sup>95</sup>]IFN- $\gamma$  (to be designated  $\gamma$ -11) is being constructed by the same general procedure. Likewise, a gene coding for [Cys<sup>95</sup>, Pro<sup>104</sup>]IFN- $\gamma$  is under construction with the threonine-specifying codon ACA (section 15 of IF-2) being replaced by the  
10 proline-specifying codon CCA.

[Glu<sup>5</sup>]IFN- $\gamma$ , to be designated  $\gamma$ -13, will result from modification of section 43 in subunit IF-3 to include the glutamate codon, GAA, rather than the aspartic acid specifying codon, GAT. Because such a change would no longer permit the presence of a BamHI recognition site at that locus, subunit IF-3 will likely need to be constructed as a composite subunit with the amino acid specifying portions of subunit IF-4, leaving no  
15 restriction site between XbaI and HindIII in the assembled gene. This analog of IFN- $\gamma$  is expected to be less acid labile than the naturally-occurring form.

The above analogs having the above-noted tryptophan and/or methionine and/or cysteine replacements are not expected to differ from naturally-occurring IFN- $\gamma$  in terms of reactivity with antibodies to the natural form or in potency of antiproliferative or immunomodulatory effect but are expected to have enhanced  
20 duration of pharmacological effects.

Still another class of analogs consists of polypeptides of a "hybrid" or "fused" type which include one or more additional amino acids at the end of the prescribed sequence. These would be expressed by DNA sequences formed by the addition, to the entire sequence coding for IFN- $\gamma$ , of another manufactured DNA sequence, e.g., one of the subunits coding for a sequence of polypeptides peculiar to LeIFN-Con, described  
25 infra. The polypeptide expressed is expected to retain at least some of the antibody reactivity of naturally-occurring IFN- $\gamma$  and to display some degree of the antibody reactivity of LeIFN. Its pharmacological activities are expected to be superior to naturally-occurring IFN- $\gamma$  both in terms of potency and duration of action.

Table VI, below, sets forth the results of studies of antiviral activity of IFN- $\gamma$  prepared according to the invention along with that of certain of the analogs tested. Relative antiviral activity was assayed in human  
30 HeLa cells infected with encephalomyocarditis virus (EMCV) per unit binding to a monoclonal antibody to IFN- $\gamma$  as determined in an immunoabsorbant assay.

35

TABLE VI

|    | <u>Interferon</u> | <u>Relative Antiviral Activity</u> |
|----|-------------------|------------------------------------|
| 40 | $\gamma$ -1       | 1.00                               |
|    | $\gamma$ -4       | 0.60                               |
|    | $\gamma$ -5       | 0.10                               |
| 45 | $\gamma$ -6       | 0.06                               |
|    | $\gamma$ -10      | 0.51                               |

50 The following example relates to modifications in the polypeptide coding region of the DNA sequences of the previous examples which serve to enhance the expression of desired products.

#### EXAMPLE 6

55

Preliminary analyses performed on the polypeptide products of microbial expression of manufactured DNA sequences coding for IFN- $\gamma$  and analogs of IFN- $\gamma$  revealed that two major proteins were produced in approximately equal quantities -- a 17K form corresponding to the complete 146 amino acid sequence and a 12K form corresponding to an interferon fragment missing about 50 amino acids of the amino terminal.

Review of codon usage in the manufactured gene revealed the likelihood that the abbreviated species was formed as a result of microbial translation initiation at the Met<sup>48</sup> residue brought about by the similarity of base sequences 3' thereto to a Shine-Delgarno ribosome binding sequence. It thus appeared that while about half of the transcribed mRNA's bound to ribosomes only at a locus prior to the initial methionine, the other half were bound at a locus prior to the Met<sup>48</sup> codon. In order to diminish the likelihood of ribosome binding internally within the polypeptide coding region, sections 33 and 34 of subunit IF-3 were re-constructed. More specifically, the GAG codon employed to specify a glutamate residue at position 41 was replaced by the alternate, GAA, codon and the CGT codon employed to specify arginine at position 45 was replaced by the alternate, CGC, codon. These changes, effected during construction of the gene specifying the  $\gamma$ -6 analog of IFN- $\gamma$ , resulted in the expression of a single predominant species of polypeptide of the appropriate length.

The following examples 7 and 8 relate to procedures of the invention for generating a manufactured gene specifying the F subtype of human leukocyte interferon ("LeuIFN-F" or "IFN- $\alpha$ F") and polypeptide analogs thereof.

#### EXAMPLE 7

The amino acid sequences for the human leukocyte interferon of the F subtype has been deduced by way of sequencing of cDNA clones. See, e.g., Goedell, et al., *Nature*, 200, pp. 20-26 (1981). The general procedures of prior Examples 1, 2 and 3 were employed in the design and assembly of a manufactured DNA sequence for use in microbial expression of IFN- $\alpha$ F in *E. coli* by means of a pBR322-derived expression vector. A general plan for the construction of three "major" subunit DNA sequences (LeuIFN-F I, LeuIFN-F II and LeuIFN-F III) and one "minor" subunit DNA sequence (LeuIFN-F IV) was evolved and is shown in Table VII below.

LeuIFN- $\gamma$  IV

**xbaI**                    -1                    1                    9                    **BstE II**  
                          Met Cys Asp Leu Pro Gln Thr His Ser Leu  
                          -43                    -41                    -42  
                          CTAGAGA ATG TGT GAT TTA CCT CAA ACT CAT TCT CTT G  
                          TCT TAC ACA CTA AAT GGA GTT TGA GTA AGA GAA CCATTG  
                          -44                    -42

5  
 10  
 15  
 20  
 25  
 30  
 35  
 40  
 45  
 50  
 55

LeuIFN-F III  
 EcoRI 10 Gly Asn Arg Arg Ala Leu Ile Leu Leu Ala Gln Met Gly Arg Ile Ser Pro Phe  
 20  
 30  
 37  
 38  
 39  
 40  
 41  
 42  
 43  
 44  
 45  
 46  
 47  
 48  
 49  
 50  
 51  
 52  
 53  
 54  
 55  
 56  
 57  
 58  
 59  
 60  
 61  
 62  
 63  
 64  
 65  
 66  
 67  
 68  
 69  
 70  
 71  
 72  
 73  
 74  
 75  
 76  
 77  
 78  
 79  
 80  
 81  
 82  
 83  
 84  
 85  
 86  
 87  
 88  
 89  
 90  
 91  
 92  
 93  
 94  
 95  
 96  
 97  
 98  
 99  
 100

BstE II  
 30 Ser Cys Leu Lys Asp Arg His Asp Phe Gly Phe Pro Gln Glu Phe Asp Gly Asn Gln Phe  
 35  
 36  
 37  
 38  
 39  
 40  
 41  
 42  
 43  
 44  
 45  
 46  
 47  
 48  
 49  
 50  
 51  
 52  
 53  
 54  
 55  
 56  
 57  
 58  
 59  
 60  
 61  
 62  
 63  
 64  
 65  
 66  
 67  
 68  
 69  
 70  
 71  
 72  
 73  
 74  
 75  
 76  
 77  
 78  
 79  
 80  
 81  
 82  
 83  
 84  
 85  
 86  
 87  
 88  
 89  
 90  
 91  
 92  
 93  
 94  
 95  
 96  
 97  
 98  
 99  
 100

Hind III  
 70 Thr Lys Asp Ser Ser Ala Thr Trp Glu Gln  
 75  
 76  
 77  
 78  
 79  
 80  
 81  
 82  
 83  
 84  
 85  
 86  
 87  
 88  
 89  
 90  
 91  
 92  
 93  
 94  
 95  
 96  
 97  
 98  
 99  
 100

40

45

50

55

30

15 \_\_\_\_\_  
 ACT CCG CTG ATG AAC G  
16 \_\_\_\_\_  
 TGA GGC GAC TAC TTG CAGCT

ECORI

114 Val Asp Ser Ile Leu Ala Val Lys Lys Tyr Phe Gln Arg Ile Thr Leu Tyr Leu 120 130

13 11 12 14

AATTCAGGC CTC GAC TCT ATT CTG GCA GTT AAA AAG TAC TTC CAG CGT ATC ACT CTG TAC CTG  
GTCCG CAG CTG AGA TAA GAC CGA CAA TTT TTC ATG AAG GTC GCA TAG TGA GAC ATG GAC

Sal I

Thr Glu Lys Lys Tyr Ser Pro Cys Ala Trp Glu Val Val Arg Ala Glu Ile Met Arg Ser Phe 140 150

9 7 5 6

ACC GAA AAG AAA TAT TCT CCG TGC GCT TGG GAA GTA GTT CGC GCT GAA ATT ATG CGT TCT TTC  
TGG CTT TTC TTT ATA AGA GGC ACG CGA ACC CTT CAT CAA GCG CGA CTT TAA TAC GCA AGT AAG

Ser Leu Ser Lys Ile Phe Gln Glu Arg Lys Glu Stp Stp 160 166 DamHI Sal I (end)

3 1 2

TCT CTG AGC AAA ATC TTC CAG GAG CGT CTG CGC CGT AAA GAA TAA TAG GATCC  
AGA GAC TCG TTT TAG AAG GTC CTC GCA GAC GCG GCA TTT CTT ATT ATC CTAGGAGCT

31

"sticky end" but, upon complementation, a Sall recognition site is not reconstituted. A full BamHI recognition site remains, however, allowing for subsequent excision of the subunit.) This first intermediate plasmid is amplified and subunit II is inserted into the amplified plasmid after again cutting with EcoRI and Sall. The second intermediate plasmid thus formed is amplified and subunit III is inserted into the amplified  
 5 plasmid cut with EcoRI and HindIII. The third intermediate plasmid thus formed is amplified. Subunit IV is ligated to an EcoRI and XbaI fragment isolated from pINT- $\gamma$ -TXb4 of Example 4 and this ligation product (having EcoRI and BstEII sticky ends) is then inserted into the amplified third intermediate plasmid cut with EcoRI and BstEII to yield the final expression vector.

The isolated product of trp promoter/operator controlled E.coli expression of the manufactured DNA  
 10 sequence of Table VII as inserted into the final expression vector was designated IFN- $\alpha$ F<sub>1</sub>.

#### EXAMPLE 8

As discussed *infra* with respect to consensus leukocyte interferon, those human leukocyte interferon  
 15 subtypes having a threonine residue at position 14 and a methionine residue at position 16 are reputed to display greater antiviral activity than those subtypes possessing Ala<sup>14</sup> and Ile<sup>16</sup> residues. An analog of human leukocyte interferon subtype F was therefore manufactured by means of microbial expression of a DNA sequence of Example 7 which had been altered to specify threonine and methionine as residues 14 and 16, respectively. More specifically, [Thr<sup>14</sup>, Met<sup>16</sup>] IFN- $\alpha$ F, designated IFN- $\alpha$ F<sub>2</sub>, was expressed in E.coli  
 20 upon transformation with a vector of Example 7 which had been cut with Sall and HindIII and into which a modified subunit II (of Table VII) was inserted. The specific modifications of subunit II involved assembly with segment 39 altered to replace the alanine-specifying codon, GCT, with a threonine-specifying ACT codon and replace the isoleucine-specifying codon, ATT, with an ATG codon. Corresponding changes in complementary bases were made in section 40 of subunit LeuIFN-FII.

The following Examples 9 and 10 relate to practice of the invention in the microbial synthesis of  
 25 consensus human leukocyte interferon polypeptides which can be designated as analogs of human leukocyte interferon subtype F.

#### EXAMPLE 9

"Consensus human leukocyte interferon" ("IFN-Con," "LeuIFN-Con") as employed herein shall mean a non-naturally-occurring polypeptide which predominantly includes those amino acid residues which are common to all naturally-occurring human leukocyte interferon subtype sequences and which includes, at  
 30 one or more of those positions wherein there is no amino acid common to all subtypes, an amino acid which predominantly occurs at that position and in no event includes any amino acid residue which is not extant in that position in at least one naturally-occurring subtype. (For purposes of this definition, subtype A is positionally aligned with other subtypes and thus reveals a "missing" amino acid at position 44.) As so defined, a consensus human leukocyte interferon will ordinarily include all known common amino acid residues of all subtypes. It will be understood that the state of knowledge concerning naturally-occurring  
 40 subtype sequences is continuously developing. New subtypes may be discovered which may destroy the "commonality" of a particular residue at a particular position. Polypeptides whose structures are predicted on the basis of a later-amended determination of commonality at one or more positions would remain within the definition because they would nonetheless predominantly include common amino acids and because those amino acids no longer held to be common would nonetheless quite likely represent the predominant  
 45 amino acid at the given positions. Failure of a polypeptide to include either a common or predominant amino acid at any given position would not remove the molecule from the definition so long as the residue at the position occurred in at least one subtype. Polypeptides lacking one or more internal or terminal residues of consensus human leukocyte interferon or including internal or terminal residues having no counterpart in any subtype would be considered analogs of human consensus leukocyte interferon.

Published predicted amino acid sequences for eight cDNA-derived human leukocyte interferon subtypes were analyzed in the context of the identities of amino acids within the sequence of 166 residues. See, generally, Goedell, et al., Nature, 290, pp. 20-26 (1981) comparing LeIFN-A through LeIFN-H and noting that only 79 amino acids appear in identical positions in all eight interferon forms and 99 amino acids appear in identical positions if the E subtype (deduced from a cDNA pseudogene) was ignored. Each of the  
 55 remaining positions was analyzed for the relative frequency of occurrence of a given amino acid and, where a given amino acid appeared at the same position in at least five of the eight forms, it was designated as the predominant amino acid for that position. A "consensus" polypeptide sequence of 166 amino acids was plotted out and compared back to the eight individual sequences, resulting in the determination that LeIFN-

F required few modifications from its "naturally-occurring" form to comply with the consensus sequence.

A program for construction of a manufactured IFN-Con DNA sequence was developed and is set out below in Table VIII. In the table, an asterisk designates the variations in IFN- $\alpha$ F needed to develop L $\alpha$ IFN-Con<sub>1</sub>, i.e., to develop the [Arg<sup>22</sup>, Ala<sup>76</sup>, Asp<sup>78</sup>, Glu<sup>79</sup>, Tyr<sup>86</sup>, Tyr<sup>80</sup>, Leu<sup>96</sup>, Thr<sup>156</sup>, Asn<sup>157</sup>, Leu<sup>158</sup>] analog of IFN- $\alpha$ F. The illustrated top strand sequence includes, wherever possible, codons noted to the subject of preferential expression in *E. coli*. The sequence also includes bases providing recognition sites for Sal, HindIII, and BstE2 at positions intermediate the sequence and for XbaI and BamHI at its ends. The latter sites are selected for use in incorporation of the sequence in a pBR322 vector, as was the case with the sequence developed for IFN- $\alpha$ F and its analogs.

10

15

20

25

30

35

40

45

50

55

### TABLE VIII

-1 1 10  
 Met-Cys-Asp-Leu-Pro-Gln-Thr-His-Ser-Leu-Gly-Asn-Arg-Arg-  
 ATG TGT GAT TTA CCT CAA ACT CAT TCT CTT GGT AAC CGT CGC  
 20 \*  
 Ala-Leu-Ile-Leu-Leu-Ala-Gln-Met-Arg-Arg-Ile-Ser-Pro-Phe-  
 GCT CTG ATT CTG CTG GCA CAG ATG CGT CGT ATT TCC CCG TTT  
 30 40  
 Ser-Cys-Leu-Lys-Asp-Arg-His-Asp-Phe-Gly-Phe-Pro-Gln-Glu-  
 AGC TGC CTG AAA GAC CGT CAC GAC TTC GGC TTT CCG CAA GAA  
 50  
 Glu-Phe-Asp-Gly-Asn-Gln-Phe-Gln-Lys-Ala-Gln-Ala-Ile-Ser-  
 GAG TTC GAT GGC AAC CAA TTC CAG AAA GCT CAG GCA ATC TCT  
 60  
 Val-Leu-His-Glu-Met-Ile-Gln-Gln-Thr-Phe-Asn-Leu-Phe-Ser-  
 GTA CTG CAC GAA ATG ATC CAA CAG ACC TTC AAC CTG TTT TCC  
 70 \* \* \* 80  
 Thr-Lys-Asp-Ser-Ser-Ala-Ala-Trp-Asp-Glu-Ser-Leu-Leu-Glu-  
 ACT AAA GAC AGC TCT GCT GCT TGG GAC GAA AGC TTG CTG GAG  
 \* \*90 \*  
 Lys-Phe-Tyr-Thr-Glu-Leu-Tyr-Gln-Gln-Leu-Asn-Asp-Leu-Glu-  
 AAG TTC TAC ACT GAA CTG TAT CAG CAG CTG AAC GAC CTG GAA  
 100 110  
 Ala-Cys-Val-Ile-Gln-Glu-Val-Gly-Val-Glu-Glu-Thr-Pro-Leu-  
 GCA TGC GTA ATC CAG GAA GTT GGT GTA GAA GAG ACT CCG CTG  
 120  
 Met-Asn-Val-Asp-Ser-Ile-Leu-Ala-Val-Lys-Lys-Tyr-Phe-Gln-  
 ATG AAC GTC GAC TCT ATT CTG GCA GTT AAA AAG TAC TTC CAG  
 130  
 Arg-Ile-Thr-Leu-Tyr-Leu-Thr-Glu-Lys-Lys-Tyr-Ser-Pro-Cys-  
 CGT ATC ACT CTG TAC CTG ACC GAA AAG AAA TAT TCT CCG TGC  
 140 150  
 Ala-Trp-Glu-Val-Val-Arg-Ala-Glu-Ile-Met-Arg-Ser-Phe-Ser-  
 GCT TGG GAA GTA GTT CGC GCT GAA ATT ATG CGT TCT TTC TCT  
 \* \* \* 160 166 Stop  
 Leu-Ser-Thr-Asn-Leu-Gln-Glu-Arg-Leu-Arg-Arg-Lys-Glu  
 CTG TCT ACT AAC CTG CAG GAG CGT CTG CGC CGT AAA GAA TAA  
 Stop  
 TAG

Table IX below sets out the specific double stranded DNA sequence for preparation 4 subunit DNA sequences for use in manufacture of IFN-Con<sub>1</sub>. Subunit LeuIFN-Con IV is a duplicate of LeuIFN-F IV of Table VIII. Segments of subunits which differ from those employed to construct the IFN- $\alpha$ F gene are designated with a "prime" (e.g., 37' and 38' are altered forms of sections 37 and 38 needed to provide arginine rather than glycine at position 22).

**TABLE IX**

LeuIFN Con IV

XbaI  
 -1 1 2 3 4 5 6 7 8 9 BstE II  
 Met Cys Asp Leu Pro Gln Thr His Ser Leu  
 CTAGAGA ATG TGT GAT TTA CCT CAA ACT CAT TCT CTT G  
 TCT TAC ACA CTA AAT GGA GTT TGA GTA AGA GAA CCATTG  
 -43 -44 -41 -42

55

50

45

40

35

30

25

20

15

10

5

LeuIFN Con III

EcoRI 10 Gly Asn Arg Arg Ala Leu Ile Leu Leu Ala Gln Met Arg Arg Ile Ser Pro Phe

39 37 38

AATTCAATTA GGT AAC CGT CGC GCT CTG ATT CTG CTG GCA CAG ATG CGT CGT ATT TCC CCG TTT  
GTAAT CCA TTG GCA GCG CGA GAC TAA GAC GAC CGT GTC TAC GCA GCA TAA AGG GGC AAA

BstE II

30 Ser Cys Leu Lys Asp Arg His Asp Phe Gly Phe Pro Gln Glu Glu Phe Asp Gly Asn Gln Phe

35 33 34

AGC TGC CTG AAA GAC CGT CAC GAC TTC GGC TTT CCG CAA GAA GAG TTC GAT GGC AAC CAA TTC  
TCG ACG GAC TTT CTG GCA GTG CTG AAG CCG AAA GGC GTT CTT CTC AAG CTA CCG TTG GTT AAG

50 Gln Lys Ala Gln Ala Ile Ser Val Leu His Glu Met Ile Gln Gln Thr Phe Asn Leu Phe Ser

60 29 30

CAG AAA GCT CAG GCA ATC TCT GTA CTG CAC GAA ATG ATC CAA CAG ACC TTC AAC CTG TTT TCC  
GTC TTT CGA GTC CGT TAG AGA CAT GAC GTG CTT TAC TAG GTT GTC TGG AAG TTG GAC AAA AGG

70 Thr Lys Asp Ser Ser Ala Ala Trp Asp Glu Hind III

79 25 26

ACT AAA GAC AGC TCT GCT GCT TGG GAC GAA  
TGA TTT CTG TCG AGA CGA CGA ACC CTG CTT TCGA

LeuIFN Con II

ECORI 80 Ser Leu Leu Glu Lys Phe Tyr Thr Glu Leu Tyr Gln Gln 90  
 23 21' 22' 24'  
 AATTCAGGCA AGC TTG CTG GAG AAG TTC TAC ACC GAG CTG TAT CAG CAG  
 GTCCGT TCG AAC GAC CTC TTC AAG ATG TGG CTC GAC ATA GTC GTC  
 Hind III

Leu Asn Asp Leu Glu Ala Cys Val Ile Gln Glu Val Gly Val Glu Glu 100  
 19' 17' 18' 20'  
 CTG AAC GAC CTG GAA GCA TGC GTA ATC CAG GAA GTT GGT GTA GAA GAG  
 GAC TTG CTG GAC CTT CGT ACG CAT TAG GTC CTT CAA CCA CAT CTT CTC

110 113 Sal I  
 Thr Pro Leu Met Asn

15 16  
 ACT CCG CTG ATG AAC G  
 TGA GGC GAC TAC TTG CAGCT

EcoRI  
 Val Asp Ser Ile Leu Ala Val Lys Lys Tyr Phe Gln Arg Ile Thr Leu Tyr Leu  
 114 120 130  
 ATTCAGGC GTC GAC TCT ATT CTG GCA GTT AAA AAG TAC TTC CAG CGT ATC ACT CTG TAC CTG  
 GTCCG CAG CTG AGA TAA GAC CGT CAA TTT TTC ATG AAG GTC GCA TAG TGA GAC ATG GAC  
 13 11 12  
 Sal I  
 Thr Glu Lys Lys Tyr Ser Pro Cys Ala Trp Glu Val Arg Ala Glu Ile Met Arg Ser Phe  
 140 150  
 ACC GAA AAG AAA TAT TCT CCG TGC GCT TGG GAA GTA GTT CGC GCT GAA ATT ATG CGT TCT TTC  
 TGG CTG TTC TTT ATA AGA GGC ACG CGA ACC CTT CAT CAA GCG CGA CTT TAA TAC GCA AGA AAG  
 9 7 5 6  
 Ser Leu Ser Thr Asn Leu Gln Glu Arg Leu Arg Arg Lys Glu Stp Stp  
 160 166  
 DamHI Sal I (end)  
 TCT CTG AGC ACT AAC CTG CAG GAG CGT CTG CGC CGT AAA GAA TAA TAG GATCC  
 AGA GAC TCG TGA TTG GAC GTC CTC GCA GAC GCG GCA TTT CTT ATT ATC CTAGGAGCT  
 3 4 2

38

the reference's summary of sequences. The four above-noted residues were retained from the native IFN- $\alpha$ F sequence to facilitate construction of subunits and assembly of subunits into an expression vector. (Note, e.g., serine was retained at position 80 to allow for construction of a HindIII site.)

Since publication of the Goedell, et al. summary of IFN- $\alpha$  subtypes, a number of additional subtypes have been ascertained. Figure 2 sets out in tabular form the deduced sequences of the 13 presently known subtypes (exclusive of those revealed by five known cDNA pseudogenes) with designations of the same IFN- $\alpha$  subtypes from different laboratories indicated parenthetically (e.g., IFN- $\alpha$ 6 and IFN- $\alpha$ K). See, e.g., Goedell, et al., *supra*; Stebbing, et al., in: *Recombinant DNA Products, Insulin, Interferons and Growth Hormones* (A. Bolon, ed.), CRC Press (1983); and Weissman, et al., *U.C.L.A. Symp.Mol.Cell Biol.*, 25, pp. 295-326 (1982). Positions where there is no common amino acid are shown in bold face. IFN- $\alpha$  subtypes are roughly grouped on the basis of amino acid residues. In seven positions (14, 16, 71, 78, 79, 83, and 160) the various subtypes show just two alternative amino acids, allowing classification of the subtypes into two subgroups (I and II) based on which of the seven positions are occupied by the same amino acid residues. Three IFN- $\alpha$  subtypes (H, F and B) cannot be classified as Group I or Group II and, in terms of distinguishing positions, they appear to be natural hybrids of both group subtypes. It has been reported that IFN- $\alpha$  subtypes of the Group I type display relatively high antiviral activity while those of Group II display relatively high antitumor activity.

IFN-Con<sub>1</sub> structure is described in the final line of the Figure. It is noteworthy that certain residues of IFN-Con<sub>1</sub> (e.g., serine at position 8) which were determined to be "common" on the basis of the Goedell, et al., sequences are now seen to be "predominant". Further, certain of the IFN-Con<sub>1</sub> residues determined to be predominant on the basis of the reference (Arg<sup>22</sup>, Asp<sup>78</sup>, Glu<sup>79</sup>, and Tyr<sup>86</sup>) are no longer so on the basis of updated information, while certain heretofore non-predominant others (Ser<sup>80</sup> and Glu<sup>83</sup>) now can be determined to be predominant.

#### EXAMPLE 10

A human consensus leukocyte interferon which differed from IFN-Con<sub>1</sub> in terms of the identity of amino acid residues at positions 14 and 16 was prepared by modification of the DNA sequence coding for IFN-Con<sub>1</sub>. More specifically, the expression vector for IFN-Con<sub>1</sub> was treated with BstEI and Hind III to delete subunit LeuIFN Con III. A modified subunit was inserted wherein the alanine-specifying codon, GCT, of sections 39 and 40 was altered to a threonine-specifying codon, ACT, and the isoleucine codon, CTG, was changed to ATG. The product of expression of the modified manufactured gene, [Thr<sup>14</sup>, Met<sup>16</sup>, Arg<sup>22</sup>, Ala<sup>76</sup>, Asp<sup>78</sup>, Glu<sup>79</sup>, Tyr<sup>86</sup>, Tyr<sup>90</sup>, Leu<sup>96</sup>, Thr<sup>156</sup>, Asn<sup>157</sup>, Leu<sup>158</sup>]IFN- $\alpha$ F, was designated IFN-Con<sub>2</sub>.

Presently being constructed is a gene for a consensus human leukocyte interferon polypeptide which will differ from IFN-Con<sub>1</sub> in terms of the identity of residues at positions 114 and 121. More specifically, the Val<sup>114</sup> and Lys<sup>121</sup> residues which duplicate IFN- $\alpha$ F subtype residues but are not predominant amino acids will be changed to the predominant Glu<sup>114</sup> and Arg<sup>121</sup> residues, respectively. Because the codon change from Val<sup>114</sup> to Arg<sup>114</sup> (e.g., GTC to GAA) will no longer allow for a Sall site at the terminal portion of subunit LeuIFN Con I (of Table IX), subunits I and II will likely need to be constructed as a single subunit. Changing the AAA, lysine, codon of sections 11 and 12 to CTG will allow for the presence of arginine at position 121. The product of microbial expression of the manufactured gene, [Arg<sup>22</sup>, Ala<sup>76</sup>, Asp<sup>78</sup>, Glu<sup>79</sup>, Tyr<sup>86</sup>, Tyr<sup>90</sup>, Leu<sup>96</sup>, Glu<sup>114</sup>, Arg<sup>121</sup>, Thr<sup>156</sup>, Asn<sup>157</sup>, Leu<sup>158</sup>] IFN- $\alpha$ F, will be designated IFN-Con<sub>3</sub>.

The following example relates to procedures for enhancing levels of expression of exogenous genes in bacterial species, especially, *E.coli*.

#### EXAMPLE 11

In the course of development of expression vectors in the above examples, the trp promoter/operator DNA sequence was employed which included a ribosome binding site ("RBS") sequence in a position just prior to the initial translation start (Met<sup>-1</sup>, ATG). An attempt was made to increase levels of expression of the various exogenous genes in *E. coli* by incorporating DNA sequences duplicative of portions of putative RBS sequences extant in genomic *E.coli* DNA sequences associated with highly expressed cellular proteins. Ribosome binding site sequences of such protein-coding genes as reported in Inokuchi, et al. *Nuc.Acids.Res.*, 10, pp. 6957-6968 (1982), Gold, et al., *Ann.Rev.Microbiol.*, 35, pp. 365-403 (1981) and Alton, et al., *Nature*, 282, pp. 864-869 (1979), were reviewed and the determination was made to employ sequences partially duplicative of those associated with the *E.coli* proteins OMP-F (outer membrane protein F), CRO and CAM (chloramphenicol transacetylase).

5'-AACCATGAGGGTAATAAATA-3'  
3'-TTGGTACTCCCATTATTTAT-5'

In order to incorporate this sequence in a position prior to the protein coding region of, e.g., the 10 manufactured gene coding for IFN-Con<sub>1</sub> or IFN- $\alpha$ F<sub>1</sub>, subunit IV of the expression vector was deleted (by cutting the vector with XbaI and BstEII) and replaced with a modified subunit IV involving altered sections 41A and 42A and the replacement of sections 43 and 44 with new segments RB1 and RB2. The construction of the modified sequence is as set out in Table X, below.

### TABLE X

20 XbaI

-1 1 2  
Met Cys Asp

RB1

|         |     |     |     |     |     |     |     |     |     |
|---------|-----|-----|-----|-----|-----|-----|-----|-----|-----|
| CTAGAAA | CCA | TGA | GGG | TAA | TAA | ATA | ATG | TGT | GAT |
| TTT     | GGT | ACT | CCC | ATT | ATT | TAT | TAC | ACA | CTA |

25 RB2

30

|     |     |     |     |     |     |     |        |
|-----|-----|-----|-----|-----|-----|-----|--------|
| 3   | 4   | 5   | 6   | 7   | 8   | 9   |        |
| Leu | Pro | Gln | Thr | His | Ser | Leu | BstEII |

41A

|     |     |     |     |     |     |     |      |
|-----|-----|-----|-----|-----|-----|-----|------|
| TTA | CCT | CAA | ACT | CAT | TCT | CTT | G    |
| AAT | GGA | GTT | TGA | GTA | AGA | GAA | CATG |

42A

Table XI, below, illustrates the entire DNA sequence in the region preceding the protein coding region of the reconstructed gene starting with the HpaI site within the trp promoter/operator (compare subunit IF-4 of Table IV).

### TABLE XI

```

5      HpaI
AAC TAG TAC GCA AGT TCA CGT AAA AAG GGT ATC TAG AAA CCA
TTG ATC ATG CGT TCA AGT GCA TTT TTC CCA TAG ATC TTT GGT

10
-1 1 2 3 4 5 6 7
Met Cys Asp Leu Pro Gln Thr His
TGA GGG TAA TAA ATA ATG TGT GAT TTA CCT CAA ACT CAT
15 ACT CCC ATT ATT TAT TAC ACA CTA AAT GGA GTT TGA GTA

8 9 BstE II
Ser Leu
20 TCT CTT G
AGA GAA CATG

```

Similar procedures were followed to incorporate sequences duplicative of RBS sequences of CRO and  
25 CAM genes, resulting in the following sequences immediately preceding the Met<sup>-1</sup> codon.

1 10 20  
CRO: GCATGTACTAAGGAGGTTGT  
CGTACATGATTCCTCCAACA

1 10 20  
CAM: CAGGAGCTAAGGAAGCTAAA  
GTCCTCGATTCCTTCGATTT

It will be noted that all the RBS sequence inserts possess substantial homology to Shine-Delgarno sequences, are rich in adenine and include sequences ordinarily providing "stop" codons.

Levels of E-coli expression of IFN-Con<sub>1</sub> were determined using trp-controlled expression vectors incorporating the three RBS inserts (in addition to the RBS sequence extant in the complete trp promoter/operator). Expression of the desired polypeptide using the OMP-F RBS duplicating sequence was at from 150-300 mg per liter of culture, representing from 10 to 20 percent of total protein. Vectors incorporating the CAM RBS duplicating sequence provided levels of expression which were about one-half that provided by the OMP-F variant. Vectors including the CRO RBS duplicating sequence yielded the desired protein at levels of about one-tenth that of the OMP-F variant.

The following example relates to antiviral activity screening of human leukocyte interferon and polypeptides provided by the preceding examples.

50 **EXAMPLE 12**

Table XII below provides the results of testing of antiviral activity in various cell lines of natural (buffy coat) interferon and isolated, microbially-expressed, polypeptides designated IFN- $\alpha$ F<sub>1</sub>, IFN- $\alpha$ F<sub>2</sub>, IFN-Con<sub>1</sub>, and IFN-Con<sub>2</sub>. Viruses used were VSV (vesicular stomatitis virus) and EMCV (encephalomyocarditis virus). Cell lines were from various mammalian sources, including human (WISH, HeLa), bovine (MDBK), mouse (MLV-6), and monkey (Vero). Antiviral activity was determined by an end-point cytopathic effect assay as described in Weck, et al., J.Gen.Virol., 57, pp. 233-237 (1981) and Campbell, et al., Can.J.Microbiol., 21, pp. 1247-1253 (1975). Data shown was normalized for antiviral activity in WISH cells.

TABLE XII

|    | <u>Virus</u> | <u>Cell<br/>Line</u> | <u>Buffy<br/>Coat</u> | <u>IFN-<br/><math>\alpha F_1</math></u> | <u>IFN-<br/><math>\alpha F_2</math></u> | <u>IFN-<br/>Con<sub>1</sub></u> | <u>IFN-<br/>Con<sub>2</sub></u> |
|----|--------------|----------------------|-----------------------|---|---|---------------------------------|---------------------------------|
| 5  | VSV          | WISH                 | 100                   | 100                                     | 100                                     | 100                             | 100                             |
|    | VSV          | HeLa                 | 400                   | 100                                     | ND*                                     | 200                             | 100                             |
| 10 | VSV          | MDBK                 | 1600                  | 33                                      | ND                                      | 200                             | 300                             |
|    | VSV          | MLV-6                | 20                    | 5                                       | ND                                      | 3                               | 20                              |
|    | VSV          | Vero                 | 10                    | 0.1                                     | ND                                      | 10                              | 0.1                             |
| 15 | EMCV         | WISH                 | 100                   | 100                                     | 100                                     | 100                             | 100                             |
|    | EMCV         | HeLa                 | 100                   | 5                                       | ND                                      | 33                              | 33                              |
|    | EMCV         | Vero                 | 100                   | 20                                      | ND                                      | 1000                            | 10                              |

\*ND - no data presently available.

It will be apparent from the above examples that the present invention provides, for the first time, an entire new genus of synthesized, biologically active proteinaceous products which products differ from naturally-occurring forms in terms of the identity and/or location of one or more amino acids and in terms of one or more biological (e.g., antibody reactivity) and pharmacological (e.g., potency or duration of effect) but which substantially retain other such properties.

Products of the present invention and/or antibodies thereto may be suitably "tagged", for example radiolabelled (e.g., with  $I^{125}$ ) conjugated with enzymes or fluorescently labelled, to provide reagent materials useful in assays and/or diagnostic test kits, for the qualitative and/or quantitative determination of the presence of such products and/or said antibodies in fluid samples. Such antibodies may be obtained from the inoculation of one or more animal species (e.g., mice rabbit, goat, human, etc.) or from monoclonal antibody sources. Any of such reagent materials may be used alone or in combination with a suitable substrate, e.g., coated on a glass or plastic particle bead.

Numerous modifications and variations in the practice of the invention are expected to occur to those skilled in the art upon consideration of the foregoing illustrative examples. Consequently, the invention should be considered as limited only to the extent reflected by the appended claims.

#### Claims

1. A method for the chemical synthesis of linear, double stranded DNA sequences of a length in excess of about 200 base pairs and coding for expression of a predetermined continuous sequence of amino acids with a selected host microorganism transformed by a selected DNA vector including the sequence, said method comprising :

(a) chemically synthesising two or more different, subunit, linear, double stranded DNA sequences of from about 100 to about 200 base pairs in length,

each different subunit DNA sequence comprising a series of nucleotide base codons coding for a different continuous portion of said predetermined sequence of amino acids to be expressed,

one terminal region of a first of said subunits comprising a portion of a base sequence which provides a recognition site for cleavage by a first restriction endonuclease, which recognition site is entirely present either once or not at all in said selected assembly vector upon insertion of the subunit therein,

one terminal region of a second of said subunits comprising a portion of a base sequence which provides a recognition site for cleavage by a second restriction endonuclease other than said first endonuclease, which recognition site is entirely present once or not at all in said selected assembly vector upon insertion of the subunit therein,

at least one-half of all remaining terminal regions of subunits comprising a portion of a

recognition site for restriction endonuclease cleavage by an endonuclease other than said first and second endonucleases, which recognition site is entirely present once and only once in said selected assembly vector after insertion of all subunits therein; and

5 (b) serially inserting each of said subunit DNA sequences prepared in step (a) into the selected assembly vector and effecting the biological amplification of the assembly vector subsequent to each insertion, thereby to form a DNA vector including the desired DNA sequence coding for the predetermined continuous amino acid sequence and wherein the desired DNA sequence assembly includes at least one unique recognition site for restriction endonuclease cleavage at an intermediate position therein.

10

2. A method according to claim 1 wherein the restriction site for endonuclease cleavage by the restriction endonuclease other than said first and second endonucleases is a palindromic six base recognition site and the desired DNA sequence assembled has at least one unique six base recognition site for restriction endonuclease cleavage at an intermediate position therein.

15

3. A method according to claim 1 wherein at least three different subunit DNA sequences are prepared in step (a) and serially inserted into said selected vector in step (b) and the desired DNA sequence obtained includes at least two unique restriction endonuclease recognition sites at intermediate positions therein.

20

4. A method according to claim 1 wherein the DNA sequence synthesised comprises an entire structural gene coding for a biologically active polypeptide.

25

5. A method according to claim 1 wherein, in the DNA sequence synthesised, the sequence of nucleotide bases includes one or more codons selected, from among alternative codons specifying the same amino acid, on the basis of preferential expression characteristics of the codon in said selected host microorganism.

30

6. A method according to Claim 2, wherein the desired DNA sequence assembled has two or more unique restriction endonuclease cleavage sites at intermediate positions therein, at least one of which has a six base palindromic recognition site.

7. A method according to Claim 4, wherein the gene codes for a human polypeptide.

35

8. A method according to Claim 4, wherein the gene codes for a polypeptide which differs from a naturally-occurring human polypeptide in terms of the identity and/or location of one or more amino acids.

40

9. A method for expressing a polypeptide, which method comprises chemically synthesising by the method of any preceding claim a DNA sequence encoding said polypeptide and causing said DNA sequence to be expressed in an organism.

#### Patentansprüche

45

1. Ein Verfahren zur chemischen Synthese von linearen, doppelsträngigen DNA-Sequenzen mit einer Länge von über etwa 200 Basenpaaren und kodierend für die Expression einer vorbestimmten kontinuierlichen Sequenz von Aminosäuren mit einem ausgewählten Wirtsmikroorganismus, der durch einen ausgewählten DNA-Vektor transformiert ist, welcher die Sequenz einschließt, wobei besagtes Verfahren umfaßt:

50

(a) chemisches Synthetisieren von zwei oder mehr unterschiedlichen, linearen, doppelsträngigen Untereinheit-DNA-Sequenzen mit von etwa 100 bis etwa 200 Basenpaaren in der Länge, wobei jede unterschiedliche Untereinheit-DNA-Sequenz eine Reihe von Nukleotidbasen-Codons umfaßt, die für einen unterschiedlichen kontinuierlichen Abschnitt besagter vorbestimmten Sequenz von Aminosäuren, die exprimiert werden soll, kodiert,

55

eine terminale Region einer ersten von besagten Untereinheiten einen Abschnitt einer Basensequenz umfaßt, der eine Erkennungsstelle für die Spaltung durch eine erste Restriktionsendonuklease bereitstellt, wobei die Erkennungsstelle vollständig entweder einmal oder überhaupt nicht in besagtem ausgewählten Assembly-Vektor bei Insertion der Untereinheit dorthinein vorhanden ist,

- eine terminale Region einer zweiten von besagten Untereinheiten einen Abschnitt einer Basensequenz umfaßt, der eine Erkennungsstelle für die Spaltung durch eine zweite Restriktionsendonuklease bereitstellt, die von besagter ersten Endonuklease verschieden ist, wobei die Erkennungsstelle vollständig einmal oder überhaupt nicht in besagtem ausgewählten Assembly-Vektor bei Insertion der Untereinheit dorthinein vorhanden ist,
- 5 wenigstens eine Hälfte aller restlichen terminalen Regionen von Untereinheiten einen Abschnitt einer Erkennungsstelle für die Restriktionsendonukleasespaltung durch eine von besagter ersten und zweiten Endonuklease verschiedene Endonuklease umfaßt, wobei die Erkennungsstelle vollständig einmal und nur einmal in besagtem ausgewählten Assembly-Vektor nach Insertion aller Untereinheiten dorthinein vorhanden ist; und
- 10 (b) der Reihe nach Inserieren jeder von besagten Untereinheit-DNA-Sequenzen, die in Schritt (a) hergestellt sind, in den ausgewählten Assembly-Vektor und Bewirken der biologischen Amplifikation des Assembly-Vektors im Anschluß an jede Insertion, um dadurch einen DNA-Vektor zu bilden, der die gewünschte DNA-Sequenz einschließt, die für die vorbestimmte kontinuierliche Aminosäuresequenz kodiert, und wobei das gewünschte DNA-Sequenz-Assembly wenigstens eine einzelne Erkennungsstelle zur Restriktionsendonukleasespaltung an einer Zwischenposition darin einschließt.
- 15
2. Ein Verfahren nach Anspruch 1, wobei die Restriktionsstelle für die Endonukleasespaltung durch die von besagter ersten und zweiten Endonuklease verschiedene Restriktionsendonuklease eine palindrome sechsbasige Erkennungsstelle ist und die gewünschte zusammengebaute DNA-Sequenz wenigstens eine einzelne sechsbasige Erkennungsstelle für die Restriktionsendonukleasespaltung an einer Zwischenposition darin aufweist.

20

  3. Ein Verfahren nach Anspruch 1, wobei wenigstens drei unterschiedliche Untereinheit-DNA-Sequenzen in Schritt (a) hergestellt werden und der Reihe nach in besagten ausgewählten Vektor in Schritt (b) inseriert werden und die gewünschte erhaltene DNA-Sequenz wenigstens zwei einzelne Restriktionsendonuklease-Erkennungsstellen an Zwischenpositionen darin einschließt.

25

  4. Ein Verfahren nach Anspruch 1, wobei die synthetisierte DNA-Sequenz ein vollständiges strukturelles Gen umfaßt, das für ein biologisch aktives Peptid kodiert.

30

  5. Ein Verfahren nach Anspruch 1, wobei, in der synthetisierten DNA-Sequenz, die Nukleotidbasensequenz ein oder mehrere Codons enthält, die unter alternativen Codons, die dieselbe Aminosäure spezifizieren, auf der Basis bevorzugter Expressionscharakteristika des Codons in besagtem ausgewählten Wirtsmikroorganismus ausgewählt sind.

35

  6. Ein Verfahren nach Anspruch 2, wobei die gewünschte zusammengebaute DNA-Sequenz zwei oder mehr einzelne Restriktionsendonukleasespaltungsstellen an Zwischenpositionen darin aufweist, von denen wenigstens eine eine sechsbasige palindrome Erkennungsstelle besitzt.

40

  7. Ein Verfahren nach Anspruch 4, wobei das Gen für ein menschliches Polypeptid kodiert.
  8. Ein Verfahren nach Anspruch 4, wobei das Gen für ein Polypeptid kodiert, das sich von einem natürlich vorkommenden menschlichen Polypeptid im Hinblick auf die Identität und/oder Stelle einer oder

45

  9. Ein Verfahren zum Exprimieren eines Polypeptids, wobei das Verfahren chemisches Synthetisieren mit dem Verfahren nach einem vorangehenden Anspruch einer DNA-Sequenz, die besagtes Polypeptid kodiert, und Bewirken, daß besagte DNA-Sequenz in einem Organismus exprimiert wird, umfaßt.

50

#### Revendications

1. Procédé pour la synthèse chimique de séquences d'ADN double brin linéaires d'une longueur dépassant environ 200 paires de bases et codant pour l'expression d'une séquence continue prédéterminée d'acides aminés avec un microorganisme hôte choisi transformé par un vecteur d'ADN choisi incluant la séquence, ledit procédé comprenant :

55

  - (a) la synthèse chimique de deux ou plus séquences d'ADN double brin linéaires, sous-unitaires, différentes, d'environ 100 à environ 200 paires de bases de longueur,

chaque séquence d'ADN sous-unitaire différente comprenant une série de codons de bases de nucléotides codant pour une portion continue différente de ladite séquence prédéterminée d'acides aminés à exprimer,

une région terminale d'une première desdites sous-unités comprenant une portion d'une séquence de bases qui fournit un site de reconnaissance pour le clivage par une première endonucléase de restriction, lequel site de reconnaissance est entièrement présent soit une fois ou pas du tout dans ledit vecteur d'assemblage choisi lors de l'insertion de la sous-unité dans celui-ci,

une région terminale d'une seconde desdites sous-unités comprenant une portion d'une séquence de bases qui fournit un site de reconnaissance pour le clivage par une seconde endonucléase de restriction autre que la première endonucléase, lequel site de reconnaissance est entièrement présent soit une fois ou pas du tout dans ledit vecteur d'assemblage choisi lors de l'insertion de la sous-unité dans celui-ci,

au moins une moitié de toutes les régions terminales restantes de sous-unités comprenant une portion d'un site de reconnaissance pour le clivage par une endonucléase de restriction autre que lesdites première et seconde endonucléases, lequel site de reconnaissance est entièrement présent une fois et seulement une fois dans ledit vecteur d'assemblage choisi après insertion de toutes les sous-unités dans celui-ci ; et

(b) l'insertion sérielle de chacune desdites séquences d'ADN sous-unitaires préparées dans l'étape (a) dans le vecteur d'assemblage choisi et l'amplification biologique du vecteur d'assemblage après chaque insertion, en formant ainsi un vecteur d'ADN incluant la séquence d'ADN souhaitée codant pour la séquence d'acides aminés continue prédéterminée et dans lequel la séquence d'ADN souhaitée, assemblée, inclut au moins un site de reconnaissance unique pour le clivage par une endonucléase de restriction en une position intermédiaire de celle-ci.

2. Procédé selon la revendication 1, dans lequel le site de restriction pour le clivage par l'endonucléase de restriction autre que lesdites première et seconde endonucléases est un site de reconnaissance à six bases palindrome, et la séquence d'ADN souhaitée, assemblée, présente au moins un site de reconnaissance unique à six bases pour le clivage par une endonucléase de restriction en une position intermédiaire de celle-ci.
3. Procédé selon la revendication 1, dans lequel au moins trois différentes séquences d'ADN sous-unitaires sont préparées dans l'étape (a) et insérées de façon sérielle dans ledit vecteur choisi dans l'étape (b), et la séquence d'ADN souhaitée, obtenue, inclut au moins deux sites de reconnaissance uniques pour le clivage par une endonucléase de restriction en des positions intermédiaires de celle-ci.
4. Procédé selon la revendication 1, dans lequel la séquence d'ADN synthétisée comprend un gène structural entier codant pour un polypeptide biologiquement actif.
5. Procédé selon la revendication 1, dans lequel, dans la séquence d'ADN synthétisée, la séquence de bases de nucléotides inclut un ou plusieurs codons choisis parmi des codons alternatifs spécifiant le même acide aminé, sur la base de caractéristiques d'expression préférentielle du codon dans ledit micro-organisme hôte choisi.
6. Procédé selon la revendication 2, dans lequel la séquence d'ADN souhaitée, assemblée, présente deux ou plusieurs sites uniques de clivage par une endonucléase de restriction en des positions intermédiaires de celle-ci, dont au moins une présente un site de reconnaissance palindrome à six bases.
7. Procédé selon la revendication 4, dans lequel le gène code pour un polypeptide humain.
8. Procédé selon la revendication 4, dans lequel le gène code pour un polypeptide qui diffère d'un polypeptide humain existant naturellement en termes de l'identité et/ou de l'emplacement d'un ou plusieurs acides aminés.
9. Procédé pour exprimer un polypeptide, lequel procédé comprend la synthèse chimique, par le procédé selon l'une quelconque des revendications précédentes, d'une séquence d'ADN codant ledit polypeptide, et l'expression de ladite séquence d'ADN dans un organisme.

## CLONING STRATEGY

FIG. 1

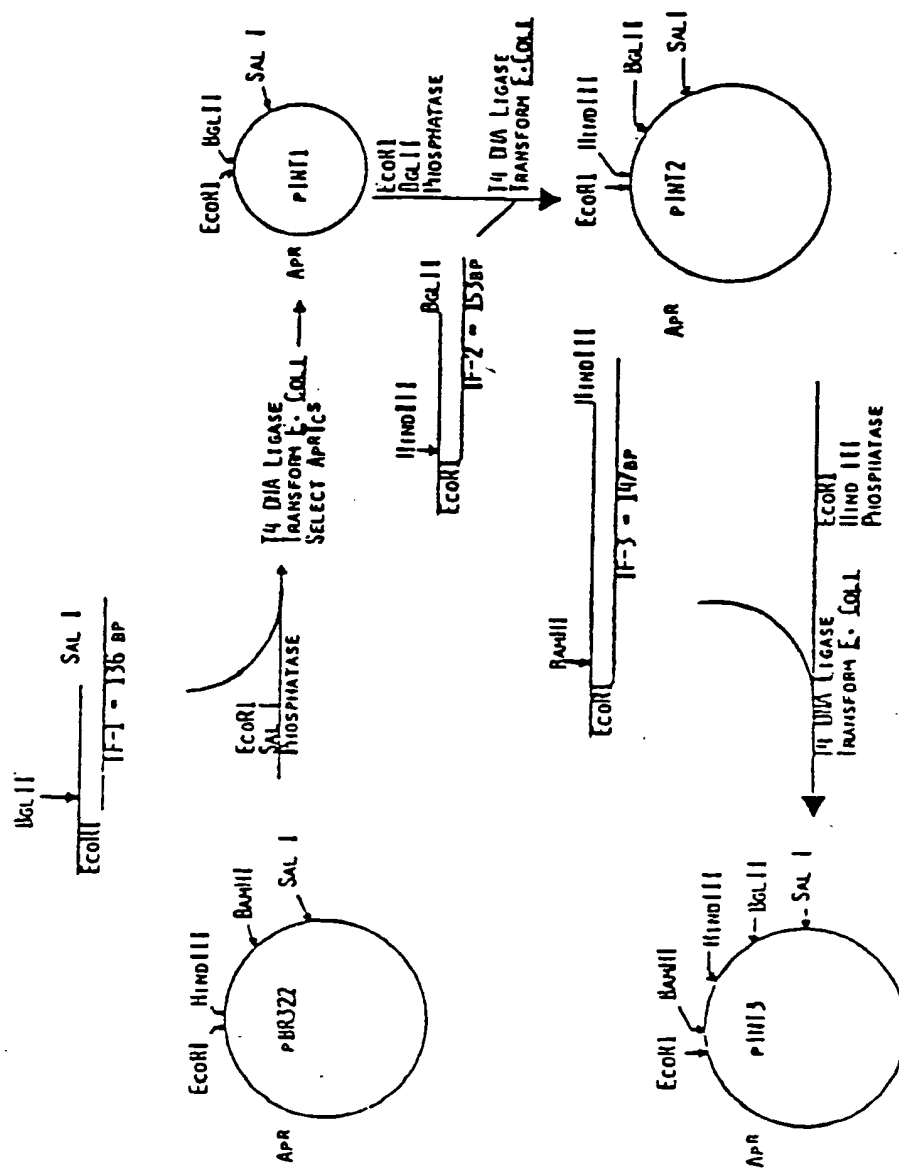


Figure 2 is a complex diagram showing a grid of data points and a series of curves. The grid is labeled with "Group" and "Subtype" on the left. The curves are labeled with "Number" and "Uniqueness" on the right. The diagram is divided into several sections, each containing a set of curves and data points. The curves are labeled with "Number" and "Uniqueness" values. The data points are labeled with "Group" and "Subtype" values. The diagram is a technical representation of a complex system, likely a biological or chemical system, showing the relationship between different components and their properties.

**FIG. 2**

